

ActivityNet Challenge 2017 Summary

Bernard Ghanem¹, Juan Carlos Niebles^{2,3}, Cees Snoek⁴, Fabian Caba Heilbron¹, Humam Alwassel¹, Ranjay Khrisna², Victor Escorcia¹, Kenji Hata², and Shyamal Buch²

¹King Abdullah University of Science and Technology

²Stanford University

³Universidad del Norte

⁴Universiteit van Amsterdam

1 Introduction

The ActivityNet Large Scale Activity Recognition Challenge is a half-day workshop held on July 26, 2017 in conjunction with CVPR 2017 in Honolulu, Hawaii. In this workshop, we want to stimulate the computer vision community to develop new algorithms and techniques that improve the state-of-the-art in human activity understanding. The data of this challenge is based on three different publicly available datasets: [ActivityNet](#), [ActivityNet Captions](#), and [Kinetics](#). The challenge focuses on recognizing high-level and goal oriented activities from user generated videos, similar to those found in internet portals.

2 Challenge Results

In this section, we list the top-3 submissions of each task along with the top-3 most innovative submissions. We also attach to this document a copy of all the papers submitted to the workshop. Please refer to the [challenge website](#) for the tasks descriptions and metrics, leaderboard, and workshop program and slides.

2.1 Task 1: Untrimmed Video Classification (ActivityNet)

Rank	Organization	Top-1 Error
1	IBUG	8.8
2	CHUK, ETHZ, and SIAT	9.8
3	Oxford Brookes University and Disney Research	18.9

Table 1: The top-3 submissions for task 1.

2.2 Task 2: Trimmed Action Recognition (Kinetics)

Rank	Organization	Average Error
1	Tsinghua and Baidu	12.4
2	CHUK, ETHZ, and SIAT	13.9
3	TwentyBN	14.4

Table 2: The top-3 submissions for task 2.

2.3 Task 3: Temporal Action Proposals (ActivityNet)

Rank	Organization	AUC
1	SJTU and Columbia	64.80
2	MSRA	64.18
3	UMD	61.56

Table 3: The top-3 submissions for task 3.

2.4 Task 4: Temporal Action Localization (ActivityNet)

Rank	Organization	Average mAP
1	SJTU and Columbia	33.40
2	CHUK, ETHZ, and SIAT	31.86
3	IC	31.82

Table 4: The top-3 submissions for task 4.

2.5 Task 5: Dense-Captioning Events in Videos (ActivityNet Captions)

Rank	Organization	Average Meteor
1	MSRA	12.84
2	U. of Science and Technology of China	9.87
3	RUC and CMU	9.61

Table 5: The top-3 submissions for task 5.

2.6 Most Innovative Submissions

Rank	Organization	Task(s)
1	BU	3 and 4
2	Tsinghua and Baidu	1 and 2
3	RUC and CMU	5

Table 6: The top-3 most innovative submissions.

Untrimmed Video Classification : submission to ActivityNet Challenge 2017

Jianqin Yin¹, Bin Wang², Xiaoli Liu¹, Bin Fang², Yanchun Wu¹,
Yanyan Bao²

jqyin@bupt.edu.cn, {wangbinth, bravebin}@tsinghua.edu.cn,{liuxlily, wuyc119,
yybaoxa}@qq.com

1. Automation School, Beijing University of Posts and Telecommunications, 100876
2. State Key Lab. of Intelligent Technology and Systems, Tsinghua University,
TNLIST, Beijing, China, 100084

Abstract

Most of traditional video action recognition methods are based on trimmed videos, which is only one action in one video. But most of videos in real world is untrimmed. In order to overcome the difficulty in some extent, we propose a method based on fusion of multiple features for untrimmed video classification task of ActivityNet challenge 2017. We use the CNN features, MBH features and stacked C3D features for classification. Then, we use one-vs-rest linear SVM to construct classifier respectively. Finally, we fuse the three results by voting to get the final recognition result.

1. Introduction

Action recognition is a hot topic in computer vision. And it is significant for computer to understand the human behavior in videos.

In order to solve this problem, we address the problem of untrimmed video classification by a fusion method based on CNN features, MBH features and C3D features.

2. System Description

2.1 Features

We use the features provided by the ActivityNet^[1] official website for untrimmed video classification.

2.1.1 Motion Boundary Histogram (MBH) Features

MBH features are extract using improved trajectories method by wang^[2]. Then the features are encoded using the GMM and Fisher Vector. To improve the system efficiency, we reduce the dimension of the MBH features with the contribution rate of

98% from 65536 to 17463.

2.1.2 Stacked C3D Features

We use the C3D^[3] features provided by ActivityNet's^[1] website. C3D features is used for extracting the temporal and spatial information of a video clip. The features were extracted every 8 frames with a temporal resolution of 16 frames. The C3D model was trained by Sport-1M dataset^[4] and it is not fine-tuned on the data of the challenge. To reduce the dimension of activations from the second fully-connected layer(fc7), PCA is used. Feature dimension is reduced from 4096 to 500. To capture the temporal information, we segment every sequence into 3 sections, and then in every section, features were mean-pooled. Then the features are stacked to obtain the final C3D features.

2.1.3 ImageNetShuffle Features

ImageNetShuffle^[5] features are extracted by using Google inception net(Google Net)^[6]. CNN features are based on the pool5 layer on two frames per second. To fuse all the frame based features, mean-pool is used across the frames followed by L1-normalization.

2.2 Classification

Based on the features description above, we apply the one-vs-rest linear SVM^[7] method to the CNN features, MBH features and C3D features respectively. And we design a fusion method to make a final decision. The fusion method can be divided into two parts. That is weight voting and hard voting. The details are as follows:

2.1 Weight Voting

If the three results are different with each other, we fuse the results by weight voting. For each type of feature, the weight for each action is learned by one-vs-rest linear SVM. For the three results, weight voting for each action respectively. And the final result is decided by the action with highest votes. Then, normalize the votes for all actions to get the confidence for each action. Finally, its confidence is decide by the max confidence.

2.2 Hard Voting

If at least two results are the same, we fuse the results by hard voting method. The final result is decide by the action with max votes. And its confidence is calculate by the following method. The motivation of the confidence calculation comes from the reliability computation.

- (1) First, for each type of feature, the confidence for each action is learn by one-vs-rest linear SVM and normalize the confidence for all actions.
- (2) Finally, calculate the confidence for the fusion result.

If one of two results are the same, set its confidence is c_1 and c_2 respectively, and the final confidence can be calculate by formula (1).

$$cof = 1 - (1 - c_1) * (1 - c_2) \quad (1)$$

If the three results are the same, set its confidence is c_1, c_2 and c_3 respectively, and the final confidence can calculate by formula (2).

$$cof = 1 - (1 - c_1) * (1 - c_2) * (1 - c_3) \quad (2)$$

3. Implementation

We use the features provided by the ActivityNet organizer (That is ImagenetShuffle features, MBH features and C3D features). Then, we apply one-vs-rest linear SVM for each type of features. Finally, we fuse the three results by the fusion method design by us.

4. Results

We report our results for the untrimmed video classification task on ActivityNet. And we evaluate our results as described in the challenge^[1]. The results is as follow:

Table 1. The results on the validation

	Validation Set	Testing set
Model	Top-1 Accuracy	Top-1 Accuracy
MBH	53.47%	-
C3D	62.83%	-
CNN	67.54%	-
MBH+C3D+CNN	74.30%	74.49%

5. Conclusion

Features from different pipelines capture different information. For example, the CNN features capture the appearance information. The MBH features and C3D features capture the dynamic information. Fusing different pipeline results can improve their performance. In this paper, we design a fusion method to improve the performance on untrimmed videos.

Reference

- [1] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 961–970, 2015. 1, 2.
- [2] H. Wang and C. Schmid. Action Recognition with Improved Trajectories. In Proc. Int. Conf. Computer Vision, pages 3551–3558, 2013.

- [3] Tran D, Bourdev L, Fergus R, et al. Learning spatiotemporal features with 3d convolutional networks, In Proceedings of the IEEE international conference on computer vision, pages 4489-4497, 2015.
- [4] Karpathy A, Toderici G, Shetty S, et al. Large-scale video classification with convolutional neural networks. In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pages 1725-1732, 2014.
- [5] Mettes P, Koelma D C, Snoek C G M. The imagenet shuffle: Reorganized pre-training for video event detection. In Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval. ACM, pages 175-182, 2016.
- [6] Szegedy C, Liu W, Jia Y, et al. Going deeper with convolutions. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 1-9, 2015.
- [7] Chih-Chung Chang and Chih-Jen Lin, LIBSVM : a library for support vector machines. ACM Transactions on Intelligent Systems and Technology, 2:27:1--27:27, 2011.

Untrimmed Video Classification: submission to ActivityNet Challenge 2017

Feng Mao
rickymf4@gmail.com

1. Model Architecture

We use CNN^[1] and C3D^[2] to extract frame level feature, and apply average pooling and LSTM to aggregate frame level feature to video level feature. Finally we adopt MoE^{[3][4]} model to do the classification.

For CNN feature, firstly, we train the inception-v1 network on the ImageNet 21k dataset. The last hidden layer before the classification layer (pool5/7x7_s1) is chosen as first part of our frame level feature. The feature vector has 1024 dimensions.

We also take advantage of the open source pre-trained C3D feature ^[5]. The 4096-dimension feature of activations from the second fully connected layer (fc7) is reduced to 500 dimensions with PCA. The video is processed with 8fps, and C3D feature is extracted configured by 16 frames per set. The C3D feature after dimension reduction forms the second part of our frame level feature.

We have two feature aggregation mechanisms. One is unsupervised average pooling. For each frame level feature ff_i^j , where i denotes the i th frame in the video, and j denotes the j th dimension of the feature. The aggregated video level feature is obtained by

$$vf^j = \sum_{i=0}^N ff_i^j / N, j \in [0, 1024]$$

For CNN feature, the pooling is overall average of all frames, for C3D, the pooling is overall average of all sets of 16 frames.

Another one is supervised aggregation. LSTM^[6] is utilized to represent the frame feature sequence. The LSTM has 2 layers with 1024 cells each. We use last layer's memory as the aggregated feature.

After extracting the video level feature, we use MoE(Mixture of experts) to model the classifier. We use 4 mixtures.

At last, we combine different features and different aggregation methods to form 4 representations:

- 1) CNN + average pooling
- 2) CNN + LSTM
- 3) C3D + average pooling
- 4) C3D + LSTM

The overall architecture is shown in Figure 1.

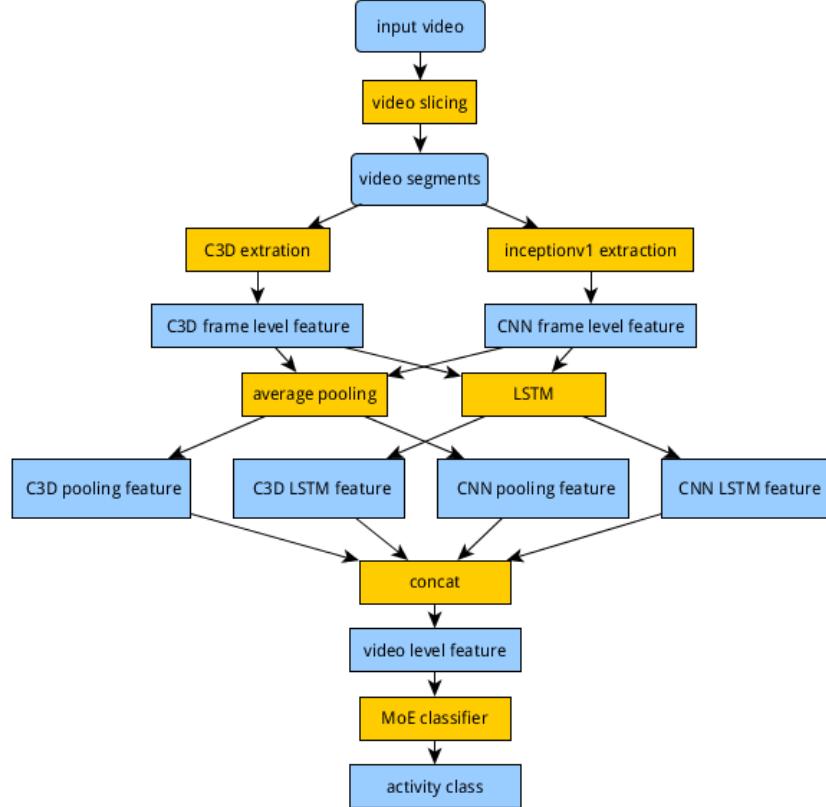


Figure 1. Model architecture

2. Training

Our training dataset is Activity Net v1.3 without external data. We utilize the pre-trained model to extract CNN and C3D frame level features for each segment from each video. The parameters needed to train are LSTM parameters for CNN and C3D and MoE parameters. Cross entropy loss between the predictions and labels is adopted:

$$\sum (z_i * -\log(f(x_i)) + (1 - z_i) * -\log(1 - f(x_i)))$$

ADAM^[7] optimizer is used with base learning rate of 0.001. The mini-batch size is 128. 1000 iterations are run before the model converges.

3. Prediction

The length of the training and validation video segments ranges from 2s to 300s, as shown in Figure 2.

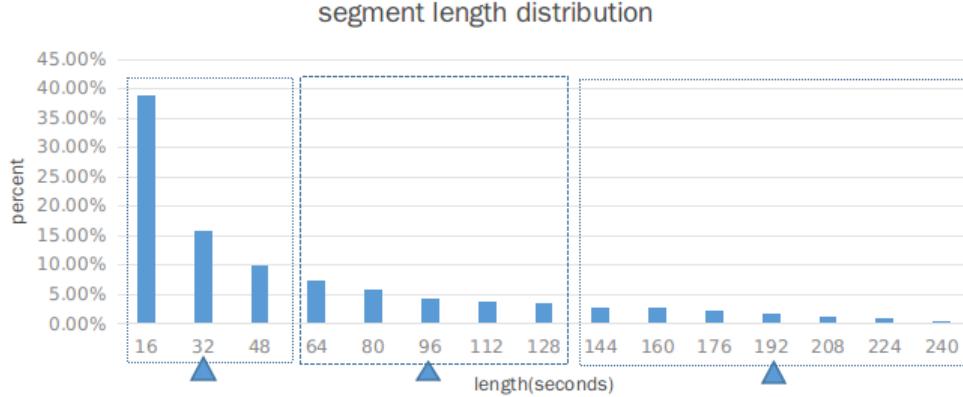


Figure 2. Segment length distribution of the training and validation set

We use 3 kinds of slide windows with size of 32, 96 192 to cover 3 ranges (3 dashed blue box in Figure 2). When inferring, we slice the video by these 3 windows without overlapping to obtain segment sets:

$$\left\{ \left\{ S_{32}^0, S_{32}^1, \dots, S_{32}^{n1} \right\}, \left\{ S_{96}^0, S_{96}^1, \dots, S_{96}^{n2} \right\}, \left\{ S_{192}^0, S_{192}^1, \dots, S_{192}^{n3} \right\} \right\}$$

Then perform prediction for all the segments, and get result sets:

$$\left\{ \left\{ f(S_{32}^0), f(S_{32}^1), \dots, f(S_{32}^{n1}) \right\}, \left\{ f(S_{96}^0), f(S_{96}^1), \dots, f(S_{96}^{n2}) \right\}, \left\{ f(S_{192}^0), f(S_{192}^1), \dots, f(S_{192}^{n3}) \right\} \right\}$$

We select top-20 classes of each f as f_{20} . The final result is the top-5 frequent classes among all the f_{20} :

$$argtop5(\bigcup_{c=1}^{200} (argtop20(f(S_i^n), c)))$$

References

- [1] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In CVPR, 2015.
- [2] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, Learning Spatiotemporal Features with 3D Convolutional Networks, ICCV 2015.
- [3] M. I. Jordan. Hierarchical mixtures of experts and the em algorithm. Neural Computation, 6, 1994.
- [4] S. Abu-El-Haija, N. Kothari, J. Lee, P. Natsev, G. Toderici, B. Varadarajan, and S. Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. arXiv preprint arXiv:1609.08675, 2016.
- [5] <http://activity-net.org/challenges/2016/download.html#c3d>
- [6] S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural Computing, 9(8):1735–1780, Nov. 1997. 2
- [7] Kingma, Diederik P. and Ba, Jimmy. Adam: A Method for Stochastic Optimization. arXiv:1412.6980 [cs.LG], December 2014.

Improve Untrimmed Video Classification by Human and Object Attention

Jiankang Deng
IBUG
Imperial College London, UK
j.deng16@imperial.ac.uk

Abstract

In this notebook paper, we introduce our submission to the untrimmed video classification task of the ActivityNet Large Scale Activity Recognition Challenge 2017. We take the method proposed by Liming Wang [12, 10] as the baseline. With data augmentation, human and object attention, and class-wise refinement, our system finally obtains the top-1 error of 8.83%.

1. Dataset and Evaluation Metric

There are 200 activity classes in the ActivityNet v1.3 dataset [1]. The training data includes 10,024 videos for training, with 15,410 activity instances. The validation set contains 4,926 videos and 7,654 activity instances. The test set has 5,044 videos with about 5000 valid download links. Only 0.1421% videos are annotated with multi-class labels, so the untrimmed video classification task is almost a single label classification problem. Besides, 75% videos have only one action, among which 78.2% actions last for more than 50% of the video length. This means more than half videos have only one action instance which lasts more than half video length. Top-1 classification error is used as the evaluation metric in the competition. Since the mean number of videos is 25 each class in the test set, every wrong video classification will decrease the performance of 4% for each class.

2. Baseline

After learning researcher Liming Wang’s talk [8], Towards efficient end-to-end architectures for action recognition and detection in videos, we decide to choose the method [12, 10] proposed by the winner of ActivityNet [1] and set up the baseline.

The convolutional networks for spatial and temporal learning are pre-trained on Image-Net classification task [2]. After testing ResNet-200 [4] and Inception V3 [7], we try ResNet-269 [13], Inception V4, Inception-ResNet-

v2 [6]. Only by visual information (appearance and motion), the combined result obtains the Top-1 accuracy of 86.9% (Top-3 accuracy 95.7%) in the validation data set. We give the easiest and hardest class names based on the Top-1 error in Table 1. The test set error reported by the server is 12.661% (submission 1). There are only a small number of classes that tend to be confused with other classes, such as Long Jump and Triple Jump, Polishing Shoes and Cleaning Shoes, Mowing the Lawn and Cutting the Grass.

3. Data Augmentation

After the performance confirmation of the baseline, we augment the training data by including the validation set. In the following experiments, there is no validation dataset to use and the evaluation result is only reported by the test server. However, the number of submission limitation is only four times. In addition, some videos from the Kinetics dataset [5] are selected into the training data. In the Kinetics dataset, there are 65 overlap classes with AcitivityNet. We further download Youtube videos by searching the keyword within the class names. We run the baseline model on these unlabeled videos. If the corresponding class is ranked within top 3, we automatically give the class label. Finally, we get a training data with balanced video numbers for each class (100 videos per class). We check the YouTube ID of the training data and confirm that the training data ID has no overlap with the test data ID.

4. Human and Object Attention

The main challenge of the untrimmed video action classification is that we do not know the accurate temporal location of the activity instance, and the unrelated video parts increase the variance of the classification problem. If some activity attention is put on the video, the classification result can be improved [9]. In [11], different activities are classified into four general types: (1) body motion only: actions fully described by human movement like “Belly dance”, (2) human object interaction [3]: actions involving specific

Easiest class	Hardest Class	Top-1 error(%)
Windsurfing	Drinking coffee	72.73
Using the pommel horse	Doing a powerbomb	54.84
Using the monkey bar	Polishing forniture	50.00
Tango	Putting on shoes	50.00
Table soccer	Removing curlers	50.00
Swinging at the playground	Rock-paper-scissors	47.06
Surfing	Gargling mouthwash	45.45
Springboard diving	Having an ice cream	42.11
Snowboarding	Polishing shoes	40.00
Snow tubing	Smoking a cigarette	36.36
Slacklining	Applying sunscreen	34.78
Skiing	Drinking beer	33.33
Shoveling snow	Washing face	33.33
Sailing	Doing nails	31.82
Rock climbing	Brushing hair	30.43
River tubing	Playing harmonica	30.43
Riding bumper cars	Painting furniture	30.00
Raking leaves	Peeling potatoes	28.57
Rafting	Cumbia	28.00
Putting in contact lenses	Cleaning shoes	26.32
Preparing pasta	Doing karate	26.32
Pole vault	Chopping wood	25.00
Volleyball	Hand washing clothes	25.00
Playing pool	Painting	25.00
Playing field hockey	Shaving legs	25.00
Playing blackjack	Using parallel bars	24.24
Playing beach volleyball	Baking cookies	24.00
Playing accordion	Playing drums	23.68
Plataform diving	Bathing dog	23.53
Plastering	Kneeling	23.53
Mixing drinks	Hopscotch	23.08
Making an omelette	Playing kickball	22.22
Longboarding	Doing crunches	21.74
Hurling	Playing saxophone	20.00
Horseback riding	Roof shingle removal	20.00
Hitting a piñata	Shot put	20.00
Hanging wallpaper	Playing flauta	19.05
Hammer throw	Swimming	18.18
Grooming dog	Preparing salad	17.86
Getting a piercing	Washing dishes	17.86
Elliptical trainer	Getting a tattoo	17.39
Drum corps	Getting a haircut	17.24
Doing motocross	Fixing bicycle	16.67
Decorating the Christmas tree	Playing guitarra	16.67
Curling	Tai chi	16.67
Croquet	Washing hands	16.67
Cleaning sink	Vacuuming floor	15.79
Clean and jerk	Waxing skis	15.79
Carving jack-o-lanterns	Doing step aerobics	15.38
Camel ride	Putting on makeup	15.38

Table 1. The easiest and hardest classes based on the Top-1 classification error. The Top-1 error of these easiest classes is zero.

objects such as “Playing piano”, (Among the 200 activity classes in ActivityNet, 76 classes are related to the ImageNet object detection classes) (3) body motion in context: body movement taking place in a specific environment like “Scuba diving”, (4) human object interaction in context: actions containing representative objects and occurring in certain context, such as “Playing beach volleyball”.

Since all of these four types are human related, we decide to improve the untrimmed activity classification by human detector. The human detector used in this competition is a GBD [13] detector trained on ImageNet and COCO dataset, with the mAP of 82.79% on ImageNet DET val2 subset. Since object related activities account for a large part of the total activity classes, we also establish a corresponding relationship between ImageNet DET dataset and ActivityNet dataset. A state-of-the-art object detection and tracking framework (detection mAP 81.7%, and tracking mAP 64.1% on the VID task of ImageNet) is utilised to get the temporal human and object tubelets. Since we do not have a state-of-the-art scene classifier, we have not use any scene context model. Based on the temporal human and object tubelets, the activity classifier is trained with the temporal human and object attention weights [9].

5. Class-wise Refinement

The top 10 accuracy of the baseline is 98.4% in the validation dataset. For each class, we fine-tune an additional two-class classifier to re-rank the videos within each class. The mean ratio of positive samples and negative samples is 1 : 9. For the fifty hardest classes, this class-wise refinement can greatly improve the top 1 accuracy.

6. Final Result

With data augmentation, human and object attention, and fifty hardest class-wise refinement, the top 1 error decreases to 9.538% (submission 2) on the test set. When acoustic features [12] are incorporated, and one hundred classes have done the class-wise refinement, the final top 1 error decreases to 8.83% (submission 3) on the test set.

7. Future Work

We plan to annotate large scale indoor and outdoor human pose from the in-the-wild activity videos. We will further investigate the relationship between articulate human pose and activities. Finally, we will do some spatial-temporal action detection experiments from untrimmed videos, and answer who/when/what at the same time.

8. Acknowledgement

Great thanks for the detailed instructions from Linchao Zhu (UTS) and Ruxing Wang (UTS). Thanks for the public

available code shared by Yuanjun Xiong (CUHK).

References

- [1] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970, 2015.
- [2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009.
- [3] G. Gkioxari, R. Girshick, P. Dollár, and K. He. Detecting and recognizing human-object interactions. *arXiv preprint arXiv:1704.07333*, 2017.
- [4] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [5] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [6] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, pages 4278–4284, 2017.
- [7] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.
- [8] L. Wang. Towards efficient end-to-end architectures for action recognition and detection in videos. 2017.
- [9] L. Wang, Y. Xiong, D. Lin, and L. Van Gool. Untrimmednets for weakly supervised action recognition and detection. *arXiv preprint arXiv:1703.03329*, 2017.
- [10] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *European Conference on Computer Vision*, pages 20–36. Springer, 2016.
- [11] Y. Wang, J. Song, L. Wang, L. Van Gool, and O. Hilliges. Two-stream sr-cnns for action recognition in videos. In *BMVC*, 2016.
- [12] Y. Xiong, L. Wang, Z. Wang, B. Zhang, H. Song, W. Li, D. Lin, Y. Qiao, L. Van Gool, and X. Tang. Cuhk & ethz & siat submission to activitynet challenge 2016. *arXiv preprint arXiv:1608.00797*, 2016.
- [13] X. Zeng, W. Ouyang, J. Yan, H. Li, T. Xiao, K. Wang, Y. Liu, Y. Zhou, B. Yang, Z. Wang, et al. Crafting gbd-net for object detection. *arXiv preprint arXiv:1610.02579*, 2016.

Untrimmed video Classification: Submission to ActivityNet

Challenge 2017

Fuchun Sun, Huaping Liu, Xuan Guo, Qi Dang, Shengsheng Wang, Menghan Guo

Department of Computer Science and Technology, Tsinghua University

{fcsun, hpliu, guoxuan, dangqi}@mail.tsinghua.edu.cn

Overview

Our method for untrimmed video classification is based on three kind of features: Improved dense trajectories (IDT)[1], C3D[2] and CNN . All features we used are provided by ActivityNet website¹¹. IDT and C3D features are preprocessed for reducing computation or normalizing the dimension of features. We use PCA to decrease the dimension of IDT features. Because of the difference of length of C3D features of each video, we simply fuse the features by average C3D features of each video clips in a video. The CNN features are extracted from pool5 layer of GoogleNet[3] on 2 frames per second. The final CNN features are fused by mean pooling.

We trained linear Support Vector Machine(SVM) classifiers and softmax classifiers in training sets for each feature.

We fuse the classification results of 6 classifiers by voting. The action which get most votes is deemed to the final result.

Reference

- [1] Wang, Heng, and Cordelia Schmid. "Action recognition with improved trajectories." *Proceedings of the IEEE international conference on computer vision*. 2013.
- [2] Tran, Du, et al. "Learning spatiotemporal features with 3d convolutional networks." *Proceedings of the IEEE international conference on computer vision*. 2015.
- [3] Szegedy, Christian, et al. "Going deeper with convolutions." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015.

¹ <http://activity-net.org/>

UTS Submission to ActivityNet Challenge 2017

Linchao Zhu Yi Yang
 University of Technology Sydney
 {zhulinchao7, yee.i.yang}@gmail.com

Abstract

We present our submission to the trimmed action recognition task in ActivityNet Challenge 2017. We mainly train the two-stream ConvNets and propose a multi-scale attention model to learn a compact video representation. We found that the multi-scale attention model outperforms average pooling on the Kinetics dataset.

1. Our Approach

We followed the basic two-stream ConvNets for video classification [3]. The ResNet-101 [1] and ResNet-152 [2] architectures are used.

In our multi-scale attention model, we introduce a memory block \mathbf{C} with shape (m, n) , where m is the number of memory slots and n is the memory size. \mathbf{C} will be updated with activations from different convolutional layers. The update procedures are as follows. We denote \mathbf{X}_i as the activations of the i th convolutional layer for the given video, \mathbf{C}_i is the center at iteration i . \mathbf{X}_i is first transformed to \mathbf{X}'_i by

$$\mathbf{X}'_i = \text{ReLU}(\text{layer_norm}(\mathbf{W}_0 \mathbf{X}_i)). \quad (1)$$

\mathbf{X}'_i is then been blended into memory \mathbf{C}_i through attention mechanism. The updated memory \mathbf{C}'_{i+1} is obtained by

$$\begin{aligned} \text{Attend}(\mathbf{X}, \mathbf{C}) &= \text{normalize}(\sum \text{softmax}(\mathbf{X}\mathbf{C}^T)(\mathbf{X} - \mathbf{C})), \\ \mathbf{C}'_{i+1} &= \text{Attend}(\mathbf{X}'_i, \mathbf{C}_i), \end{aligned} \quad (2)$$

where the normalization function is SSR normalization followed intra normalization and ℓ_2 normalization. \mathbf{C}'_{i+1} is then transformed to the next memory \mathbf{C}_{i+1} with

$$\mathbf{C}_{i+1} = \text{ReLU}(\text{layer_norm}(\mathbf{W}_1 \mathbf{C}'_{i+1})). \quad (3)$$

After iterations N , \mathbf{C}_N is flattened and used for classification.

2. Experiments

To train the RGB network, we initialize the weights from the pre-trained ImageNet models. The flow net is initialized

Model	RGB	Flow	RGB+Flow
ResNet V1 101	69.9 / 88.2	60.5 / 81.1	72.6 / 89.8
ResNet V2 152	70.2 / 88.5	—	—
Multi-scale Attention	71.1 / 90.0	—	—

Table 1. Single checkpoint, single scale performance.

Model	Validation	Test
RGB models	72.4 / 90.4	—
RGB+Flow	74.2 / 91.0	—

Table 2.

with the trained RGB model [4]. We used SGD with momentum 0.9, and the initial learning rate is 0.01. The batch size is set to 256 for ResNet-101 and 128 for ResNet-152. The learning rate decays 0.1 every 200,000 iterations.

To train the multi-scale attention network, we randomly sample 16 frames from a video and $1/k$ positions are sampled from the feature map with size (k, k) . We used 3 feature maps from ResNet-101, which are the activations from block 2, block 3, and block 4. The results are shown in Table 1.

We fused all RGB models by average fusion. The RGB scores and flow scores are also averaged to obtain the final scores. The fusion results are shown in Table 2.

References

- [1] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [2] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016.
- [3] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014.
- [4] L. Wang, Y. Xiong, Z. Wang, and Y. Qiao. Towards good practices for very deep two-stream convnets. *arXiv preprint arXiv:1507.02159*, 2015.

NUS Submission to ActivityNet Challenge 2017: Evaluating frame-based and spatio-temporal CNN features for video classification

An Tran, Loong-Fah Cheong

Department of Electrical & Computer Engineering, National University of Singapore

an.tran@u.nus.edu, eleclf@nus.edu.sg

Abstract

This paper presents our findings when participating in ActivityNet challenges in untrimmed and trimmed video classification on ActivityNet and Kinetics dataset. Flow data are not as reliable as RGB data in challenging datasets such as ActivityNet. Furthermore, frame-based BN-Inception architecture performs better than spatio-temporal based C3D models.

1. Introduction

With the appearance of big video datasets (such as ActivityNet [1], Kinetics [4]), it would be interesting to evaluating different convolutional neural networks (CNNs) for action recognition. Two popular deep learning based representations for videos are frame-based CNNs features (*e.g.*, Inception [2]) and spatio-temporal CNNs features (*e.g.*, C3D [7]). In recent work, Wang *et al.* [8] shows that frame-based CNNs features achieves the state-of-the-art results on small video datasets such as UCF101 [6], HMDB51 [5]. In this submission to the challenge, we aims to evaluate the performance of frame-based CNNs and spatio-temporal CNNs on large video dataset.

2. Our approach

Models. For frame-based features, we deploy BN-Inception [2] architecture with temporal-segment-networks (TSNs) proposed in [8]. For spatio-temporal features, we utilize the C3D architecture [7]. In C3D architecture, all convolution and pooling layers are made up of 3D operations [7]. C3D network has 8 convolution and 3 fully connected layers. C3D architecture is then given by: $C(3, 64, 1) - RL - P(2, 2) - C(3, 128, 1) - RL - P(2, 2) - C(3, 256, 1) - RL - C(3, 256, 1) - RL - P(2, 2) - C(3, 512, 1) - RL - C(3, 512, 1) - RL - P(2, 2) - C(3, 512, 1) - RL - C(3, 512, 1) - RL - FC(4096) - D(0.9) - FC(4096) - D(0.8) - FC(101)$. Table 1 shows the number of parameters in BN-Inception

and C3D architecture. As can be observed, C3D has more parameters than BN-Inception.

We also reports performances of some long-term temporal convolutions (LTC) models. LTC models develop the ideas of C3D architectures into longer temporal dimension by reducing spatial dimension and increasing the temporal length. Inception model is initialized from pre-trained model on ImageNet and C3D weights are initialized from 1M-Sports dataset.

Regularization. For regularizing, the TSNs [8] utilize batch normalization and high dropout ratios. We also set high dropout ratios for C3D networks with ratios of 0.9 and 0.8 for fully connected layers fc6 and fc7 respectively.

Input. The inputs for TSN models are 1 RGB frame for a segment of spatial-CNN stream and 5 stacked optical flows frames for a segment of temporal-CNN models. The TSN networks average 3 segments of frames. Hence, TSN networks operates on 3 RGB frames and 15 flows frames. The C3D networks operates on 16 continuous frames for both RGB and optical flows. The flows are extracted from OpenCV implementations of TV-L1 [9]. All the RGB and flow images are saved in resolution (128, 171). The flows have been compensated to remove camera motions. For the BN-Inception architecture in TSNs, the input images are resized into (256, 340) on the fly in our modifications of Caffe software [3]. Although it is convenient for working with both BN-Inception and C3D architecture, it might slightly reduce performance of BN-Inception architecture because of lower quality resized image frames. Spatial and temporal models are trained individually.

Data augmentation. Data augmentation is shown to help deep convolutional models prevent severe over-fitting. For training, we adopt widely used data augmentation such as corner cropping, horizontal flipping, scale-jittering [8].

3. Results

Table 2 shows performance of different models on ActivityNet validation set. The results show that both C3D and TSN-BN-Inception obtain better performances on RGB

	BN-Inception [2]	C3D [7]
parameters	10,373,765	78,409,573

Table 1. Number of parameters of different convolution networks: CaffeNet, and C3D.

data than on compensated optical flows. For examples, rgb-C3D-size112-len16 out-performs flow counterpart flow-C3D-size-112-len16 17.09% (65.36% vs. 48.27% top-1 accuracy). The similar phenomenon happens with TSN-BN-Inception model (72.73% vs. 53.40%). On contrary to performances of deep models [8] on UCF101 and HMDB51, in challenging datasets such as ActivityNet, the compensated flows are not reliable as RGB data. It is because motion patterns in ActivityNet and Kinetics dataset are more complex than small datasets (*e.g.*, UCF101, HMDB51).

The second observation is that Inception models generally outperform C3D models both in RGB modality, while LTC-C3D models have more advantages in flow modality with longer temporal length. Furthermore, two-stream BN-Inception model perform better two-stream C3D. It indicates it would be hard to fit C3D models with spatio-temporal data.

The third observation is that more information you can feed into a C3D models, more successful the model is. The C3D model has better performance if we increase temporal length both RGB and flow data. As can be observed, spatial resolution is more important for rgb-C3D (spatial-stream), while temporal resolution is more important for flow-C3D (temporal-stream). It can be explained as spatial-stream exploits context in whole videos, while temporal-stream focus more long-term human motion informations (*e.g.*, human silhouettes, or shapes).

Due to limited time and resources, we take a chance to submit our evaluations on test set into the challenge server with only rgb-LTC-maxpool-size112-len32 model. We got about 34% top-1 on ActivityNet untrimmed video classification task, and 30% average error (top-1 and top-5 errors) on Kinetics trimmed video classification.

4. Conclusions

This paper reports some comparisons between current frame-based BN-Inception and spatio-temporal C3D models. Without regards of models, RGB data are more reliable than (compensated) optical flows in challenging datasets such as ActivityNet and Kinetics. With regards of models, frame-based BN-Inception currently perform better than spatio-temporal based C3D models.

References

- [1] B. G. Fabian Caba Heilbron Victor Escorcia and J. C. Niebles. ActivityNet: A Large-Scale Video Benchmark for Human Ac-

Models	Top-1 acc. (%)	Top-3 acc. (%)
rgb-C3D-size112-len16	65.36	81.43
rgb-LTC-maxpool-size112-len32	67.29	82.22
rgb-LTC-maxpool-size58-len32	56.13	72.80
flow-C3D-size-112-len16	48.27	65.96
flow-LTC-maxpool-size112-len32	59.10	74.53
flow-LTC-maxpool-size58-len32	53.65	70.19
rgb-TSN-BN-Inception	72.73	87.16
flow-TSN-BN-Inception	53.40	70.60
Two-stream C3D	70.36	83.95
Two-stream TSN-BN-Inception	73.25	86.30

Table 2. Performance of TSNs and C3D models on ActivityNet untrimmed video validation set. We show the results in top-1 and top-3 accuracy metrics. The suffix size112 means that the C3D model has input image of size (112, 112) and len16 means that the C3D model has temporal length of 16.

- tivity Understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
- [2] S. Ioffe and C. Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. feb 2015.
 - [3] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional Architecture for Fast Feature Embedding. *arXiv preprint arXiv:1408.5093*, jun 2014.
 - [4] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman. The Kinetics Human Action Video Dataset. may 2017.
 - [5] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: A large video database for human motion recognition. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2556–2563, 2011.
 - [6] K. Soomro, A. R. Zamir, and M. Shah. UCF101 : A Dataset of 101 Human Actions Classes From Videos in The Wild. Technical Report November, 2012.
 - [7] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning Spatiotemporal Features with 3D Convolutional Networks. In *The IEEE International Conference on Computer Vision (ICCV)*, dec 2015.
 - [8] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal Segment Networks: Towards Good Practices for Deep Action Recognition. In *ECCV 2016 - European Conference on Computer Vision*, aug 2016.
 - [9] C. Zach, T. Pock, and H. Bischof. A duality based approach for realtime tv-l1 optical flow. In *In Ann. Symp. German Association Patt. Recogn*, pages 214–223, 2007.

NIAD Submission to ActivityNet Challenge 2017

Kensho Hara, Hirokatsu Kataoka, Yusuke Goutsu, Yutaka Satoh
National Institute of Advanced Industrial Science and Technology (AIST)
Tsukuba, Ibaraki, Japan
`{kensho.hara, hirokatsu.kataoka}@aist.go.jp`

Ryozo Yamashita
DENSO CORPORATION
Chuo-ku, Tokyo, Japan
`RYOZO_YAMASHITA@denso.co.jp`

Abstract

In this paper, we introduce our method for ActivityNet Challenge 2017. We mainly focus on the temporal action proposals task. We use a 3D convolutional neural network (CNN) to generate action proposals. The CNN outputs actionness scores, begin times and duration of action proposals. We experimentally evaluate the performance of our method.

1. Introduction

We mainly focus on the temporal action proposals task in ActivityNet Challenge 2017. We use a 3D convolutional neural network (CNN) [1] to generate action proposals. The inputs of the network are short video clips of multiple temporal scales, and the CNN outputs an actionness score, begin time and duration of an action proposal. We generate action proposals by applying the CNN to entire video clips.

We submitted our method to task 1: untrimmed video classification, task 2: trimmed action recognition, and task 3: temporal action proposals.

2. Methodology

Our method is based on 3D CNNs. 3D CNNs apply 3D convolutional filters to videos in the spatio-temporal 3D space to extract spatio-temporal features of videos. 3D CNNs recently achieve good performances for activity recognition by using large scale datasets such as Kinetics [2].

The network architecture of our method is shown in Fig. 1. The inputs of the network are short video clips of multiple temporal scales, Following to the last layer, the network has the layer that concatenates activations of the last



Figure 1. Network architecture of our method for the temporal action proposals task.

layer based on each scale to combine the multi-scale information. The concatenated layer connects to the output layer that includes actionness score, begin time and duration of an action proposal.

Whereas our method requires the multi-scale inputs, CNNs require fixed-size inputs. In addition, convolutions for large scale inputs are computationally expensive. Therefore, we thin frames at regular intervals to align the temporal size of multi-scale inputs. This alignment enables CNNs to process multi-scale inputs with small computational costs.

In the training step, we randomly generate positive and negative samples. Each sample has multiple scales. The positive samples overlap with activities larger than 0.7, and the negative samples do not overlap with activities larger than 0.3. The overlap ratio is calculated using the smallest scale. The number of negative samples is the same as that of the positive samples. In addition to a class label, a begin time and duration of corresponding action are annotated to a positive sample. The annotation of begin time is relative offset from the temporal center of sample to the begin time of action. The offset is normalized by the duration of action. The annotation of duration is log-space shift relative to the smallest scale. We minimize an multi-task loss function using SGD with momentum to train the CNN. The loss function includes a classification loss (cross-entropy) and

Table 1. Results of our method.

Method	Untrimmed		Trimmed		Proposals	
	val	test	val	test	val	test
<i>C3D</i>	61.0	-	-	-	51.6	51.7
<i>ResNet</i>	64.8	65.3	-	68.9	-	-

localization loss (smooth L1), similar with Faster R-CNN [3].

We adopt the sliding window manner to generate action proposals in the test step, similar with [4]. Our method outputs score, begin time, and duration of a proposal at each position. We also perform thresholding for the scores and non-maximum suppression.

3. Experiments

In this section, we first describe the setting of our method, and then show experimental results of our method. We used the CNNs that are based on C3D [1] and ResNet-34 [5]. The C3D and ResNet-based models are pre-trained using the Sports-1M [6] and Kinetics [2] datasets, respectively. We changed the 2D convolution kernels to 3D ones of the ResNet-based model. The spatial scales of inputs are 112×112 . We cropped each frame at the center position and resized it to the scale. The temporal scales of inputs are $\{16, 128, 1024, 8192\}$. The larger scales are thinned to 16 frames when inputting to the CNN.

The models for untrimmed and trimmed action recognition tasks have a softmax layer for action recognition instead of the concatenation layer. For the untrimmed task, we fine-tuned the CNN models on the ActivityNet dataset. To recognize actions of each video, we input non-overlap 16 frame clips to the model, and the final recognition score is obtained by averaging the softmax scores for the clips.

Table 1 shows the results of three tasks. The results of untrimmed, trimmed, and proposals tasks are top-1 accuracy, average accuracy of top-1 and top-5, and area under average recall vs. average number of proposals per video curves, respectively.

4. Conclusion

In this paper, we describe the 3D CNN based method for the temporal action proposals task. The CNN outputs actionness scores, begin time and duration of action proposals using temporal multi-scale clips. Our method is applied to entire videos by the sliding window manner.

References

- [1] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, “Learning spatiotemporal features with 3d convolutional networks,” in *Proceedings of the International Conference on Computer Vision (ICCV)*, pp. 4489–4497, 2015.
- [2] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman, “The kinetics human action video dataset,” *arXiv preprint*, vol. arXiv:1705.06950, 2017.
- [3] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” in *Advances in Neural Information Processing Systems 28*, pp. 91–99, 2015.
- [4] Z. Shou, D. Wang, and S.-F. Chang, “Temporal action localization in untrimmed videos via multi-stage CNNs,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1049–1058, 2016.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [6] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, “Large-scale video classification with convolutional neural networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1725–1732, 2014.

CUHK & ETHZ & SIAT Submission to ActivityNet Challenge 2017

Yue Zhao¹, Bowen Zhang², Zhirong Wu¹, Shuo Yang¹, Lei Zhou², Sijie Yan¹, Limin Wang³, Yuanjun Xiong¹, Dahua Lin¹, Yu Qiao², and Xiaoou Tang¹

¹Department of Information Engineering, The Chinese University of Hong Kong

²Shenzhen Institutes of Advanced Technology, CAS, China

³Computer Vision Laboratory, ETH Zurich, Switzerland

Abstract

This paper presents the method for our submission to the tasks of trimmed video action recognition, untrimmed video classification, and temporal action localization of ActivityNet Challenge 2017. In the trimmed action recognition task, we integrate short-term temporal information with 3D models and long-term temporal information with temporal segment networks [18]. We also explore incorporating multi-modal cues such as acoustic signal and body pose. We find that models trained on Kinetics can serve as a better pretrained model than ImageNet pretrained models, and greatly boost the performance for other video understanding tasks. For temporal action localization, we benchmarked our recently proposed structural segment networks [21]. For untrimmed video classification, we study different temporal aggregation schemes such as UntrimmedNet [17].

1. Introduction

Human action understanding including recognition and detection has drawn a significant amount of attention from the academic community [12]. Publicly available benchmarks such as ActivityNet [3], THUMOS [7], and UCF-101 [13] have contributed to spurring interest and progress in action recognition and detection research. As a result, a number of deep-learning based frameworks, including TSN [18], C3D [16], SSN [21] have been proposed and achieved impressive results on the benchmark datasets. The success of these methods cannot do without the highly discriminative feature representation generated by convolutional network in a data-driven manner.

To obtain robust representation, large number of data and effective learning strategies are required. Previous methods transfer image-level object recognition representation

to video-level action recognition representation. A more natural way is to learn general video-level recognition representation and transfer these representation to other video analysis applications. Recently, a large scale action recognition dataset named Kinetics [9] is released which contains 300,000 trimmed video with 400 action categories. This dataset can serve as a good source for video-level representation learning.

In this report, we focus on learning video-based representation by conducting action recognition using Kinetics dataset. We show that these features can be easily transferred to other video analysis tasks, i.e. action detection and other datasets, Activity-Net [3]. Different from image-level recognition, temporal structure is essential for action recognition. Temporal segment networks [18] have been proven to be an effective learning strategies to model temporal structure in actions. We use temporal segment networks [18] as our base model and introduce various techniques to enhance feature learning e.g. re-balancing the sampling when training, (un)freezing BN units and tuning the magnitude of dropout depending on the size of dataset. Apart from effective training strategies, new network units, including 3D pooling and 3D convolution, are also explored to model the short-term temporal information.

While video understanding certainly depends on the visual analysis, we notice that other sources of information, such as acoustic signal and human pose that come along provide complementary information. To exploit such information, we develop a deep network called Pose CNN to infer human activity from the coordinates of body's key-point, and another one called Audio CNN to derive complementary features from the spectrograms. Combining all of the visual, pose, and acoustic models, we attain a high recognition accuracy (top-1 error of 9.79% on Activity-Net Untrimmed and average error of 13.93% on Kinetics testing set).

These features learned from Kinetics [9] can be easily

transfer to other video tasks, e.g., temporal action detection. We follow the framework of SSN [21] to conduct temporal action detection, achieving superior detection results (average mAP of 31.86%) over previous ImageNet-pretrained models.

The remaining part of this report is organized as follows. Section 2, 3, 4 present our approach in three tasks of ActivityNet 2017 challenge, namely, trimmed action recognition (Task 2), untrimmed video classification (Task 1), and temporal action localization (Task 4). Section 5 concludes this work.

2. Trimmed Action Recognition

In this year’s challenge, the trimmed video classification task in conducted on the Kinetics dataset [9]. There are around 300,000 temporally trimmed videos used for this task. Each video lasts around 10 seconds. Our efforts on this task are focused on how to utilize temporal information, both short-term and long-term, to learn better features for video-based action recognition. Now we describe how we devise our approach based on this principle.

2.1. Long-term Temporal Information

We follow the pipeline of temporal segment networks (TSN) [18] to model the long-term temporal information. Using TSN, we divide every 10-second video to a fixed number of segments. During training, one short snippet is sampled from each segment, which forms a sparse snippet sampling scheme. The snippet-wise prediction is then aggregated using different strategies such as average pooling or top-k pooling. During testing, we follow the standard procedure of using 25 frames uniformly extracted from the testing videos and average the predictions.

Since videos of different classes are unevenly distributed in the training set but balanced in the validation and testing set, we propose to re-balance data when training. For the Inception BN [8]-based flow stream, the top-1 accuracy is raised from 60.55% to 62.07%.

We experiment with several recent network architectures such as Inception V3 [15], ResNet [6], and Inception-ResNet v2 [14]. Besides, we also compare the models with ImageNet pretraining and the models trained from scratch on Kinetics [9]. We find that models pre-trained from ImageNet [2] achieve higher accuracy than models trained from scratch. We also experiment with models pre-trained from WebVision [10] and achieved on-par results as ImageNet pretraining. This suggests that pretraining on image dataset is still beneficial even at the current scale of Kinetics dataset.

2.2. Short-term Temporal Information

In the original TSN [18], the underlying ConvNet models are using 2D inputs, *i.e.*, RGB frames and optical flow

Table 1. Performance comparison of different model architectures on Kinetics [9].

Models	RGB		Flow	
	Top-1	Top-5	Top-1	Top-5
Inception BN	69.09%	88.67%	62.93%	83.67%
Inception V3	72.51%	90.22%	63.07%	84.52%
ResNet-200	70.73%	89.19%	-	-
IR-10	72.26%	90.52%	-	-

images. To further exploit short-term temporal variation between neighboring frames, we propose to modify the original 2D model into a 3D version. The number of input frames in each snippet is changed to 8 or 32 instead of 1. Two types of 3D models are experimented: (1) Pooling layers are inflated to be 3-dimensional while convolutional layers remain 2-dimensional. (2) Both pooling and convolutional layers are inflated to be 3-dimensional. We found that using 3D pooling alone with TSN can get a good boost of performance compared with 2D models. The results are illustrated in Table 3.

Table 2. Performance of different spatial-stream 3D-pool architectures on Kinetics [9].

TSN segments	base model	Top-1	Top-5
1 (no TSN)	Inception BN	69.97%	88.90%
2 (TSN)	Inception BN	71.12%	89.70%
TSN	Inception V3	73.83%	91.04%

Table 3. Performance of different spatial-stream 3D-conv architectures based on Inception BN [8].

# of frames	Top-1	Top-5
8	69.64%	88.95%
32	65.68%	86.73%

2.3. Pose Information

Besides pixel level information, we are also interested in how human poses function in action recognition at this large scale. To this end, we make a preliminary attempt of using pose estimation output for action recognition on Kinetics [9]. We use OpenPose [1] as our pose estimator to extract 18-key-point body pose features $(x, y, s)_{i \in \{1, \dots, 18\}}$ from every 5-th frame of the videos. This results to around 60 frames of data for each video. We design a 1D ConvNet model that takes in processed joints coordinates and outputs action categories. The model architecture is reflected in table 4.

In training, we sample 48 continuous frames from one video as augmentation and feed the normalized feature into a 1-d CNN to produce prediction. For evaluation, all frames are fed into the network in convolutional manner and we pool the output of the final classification layer to obtain the action prediction.

Table 4. Model architecture for the Pose CNN.

layer name	output size	(full pose)
conv1	48x1	3x1, 128
		3x1, 256
conv2	48x1	concat with parts
		3x1, 256
		3x1, 256
conv3	24x1	3x1, 256, stride 2
		3x1, 256
conv4	12x1	3x1, 256, stride 2
		3x1, 256
		3x1, 512, stride 2
conv5	6x1	3x1, 1024
		3x1, 1024
		3x1, 1024
avg. pooling	1x1	6x1
dropout		
fc	400x1	
global pool		
softmax		

Table 5. Performance of different methods for pose on Kinetics [9].

method	Top-1
pose + encoding + SVM	13.96%
pose + 1d CNN	20.13%
hierarchical pose + 1d CNN	24.11%

To our relieve, even with the ultra low dimensionality of the joints data, our models are still able to achieve more than 20% top-1 accuracy and around ??% top-1 accuracy. We also find it obtains better results if we at first treat pose feature of different parts separately and then aggregate them in subsequent layers. We collect pose feature from 5 parts, namely upper body, lower body, head, arms, limbs. The feature map are concatenated after two convolutional units. This hierarchical design achieves an accuracy 24.11%. Further fusion shows that the pose-based prediction is complementary to that of the two-stream networks. The results of using pose information are summarized in Table 5.

2.4. Acoustic Information

Audio signals in a video carry important cues for recognizing certain types of activity. We use audio-based CNN similar to the one we used in last year’s challenge [20] to extract acoustic cues to facilitate action recognition. The raw 44100-Hz mono audio signal is transformed to spectrograms in log-scale and converted to grayscale time-frequency feature maps. The network architecture is Inception BN [8], whose weights are initialized by the same technique used on the temporal networks[18]. We also experiment with TSN-style training by sampling a spectrogram patch within 3 snippets and predicting labels based on con-

sensus. The results of acoustic models are summarized in Table 6. Although the accuracy of acoustic models seem inferior, we do find them to form important components in a ensemble. This again confirms that audio information is complementary to visual information for action recognition.

Table 6. Performance of audio CNN on Kinetics [9].

TSN segments	Top-1	Top-5
1 (no TSN)	19.44%	36.48%
3 (TSN)	23.10%	41.06%

Table 7. Results on Kinetics [9] test set. The avg. error is the average of top-1 and top-k error.

Models	avg. error
model ensemble	0.1393

3. Untrimmed Video Classification

The second task we participated in is the untrimmed video classification task which aims at classifying temporal untrimmed videos from 200 categories in the ActivityNet [3] dataset. On this task we focus on examining the benefit of feature learning on the large scale trimmed video dataset, such as Kinetics [9]. Thus we follow the TSN [17] framework and adapt it to training on annotated action instances. For testing we adopt the aggregation scheme introduced by [19]. The outcome of experiments is encouraging. Models trained from Kinetics [9] can be directly adapted here and render substantial improvement over ImageNet [2] pre-trained models, shown in 8.

Table 8. Performance comparison of different model architectures on ActivityNet v1.3 [3] val set.

Models	source	Top-1 Accuracy	
		RGB	Flow
Inception BN	ImageNet	75.44%	63.66%
	Kinetics	81.95%	73.45%
Inception V3	ImageNet	78.81%	64.68%
	Kinetics	83.46%	75.58%
ResNet 101	ImageNet	80.71%	63.80%
	Kinetics	83.72%	75.90%

Likewise in the trimmed action recognition task, we also adopt the 3D input based models for this task. These models are pretrained on the Kinetics dataset and fine-tuned on the untrimmed video classification task.

3.1. UntrimmedNets

Understanding that temporal annotation may be labour intensive to obtain, we experimented with the recently proposed UntrimmedNet framework [17] to directly learn classification models from the whole untrimmed videos. This

Table 9. Performance of 3D models on ActivityNet v1.3 [3] val set. All are initialized with Kinetics-pretrained models.

Model	# of frames	Top-1	Top-3
3D Conv + 3D Pool	8	79.50%	92.51%
3D Conv + 3D Pool	32	81.81%	92.94%
2D Conv + 3D Pool	8	82.14%	93.58%

framework features a carefully designed selection mechanism which embedded into the end-to-end training process to find and learn from useful video segments. We find that it can achieve comparable results (see Table 10) as those models trained with temporal annotations. Additionally, we found that they are complementary with those trimmed training models. These findings suggest that (1) we may be able to learn action recognition from large numbers of untrimmed videos directly

Table 10. Performance of UntrimmedNet on ActivityNet v1.3 [3] val set.

Models	Top-1	Top-3
Inception BN	81.91%	93.40%
Inception V3	82.76%	94.15%

Table 11. Results on ActivityNet v1.3 [3] test set.

Models	Top-1
Ensemble of models pre-trained on ImageNet	88.28%
Ensemble of models pre-trained on Kinetics	90.21%

4. Temporal Action Localization

4.1. Structure Segment Networks

We use structure segment networks, described in [21] in this task. The model adopts the "proposal + classification" scheme in the modern object detection architecture like Fast R-CNN [4]. To generate temporal proposals, we first train a TSN-based binary classifier. The prediction output can be viewed as a series of actionness probabilities, upon which a set of temporal proposals can be generated using temporal actionness grouping [21].

During the training phase, video frames and TAG proposals are fed into SSN. For each proposal, we first expand the boundary by extending forward and backward. The augmented proposal is thus composed of three stages. Extracted similar to TSN, snippet-level features are aggregated via structured temporal pyramid pooling. We divide proposal samples into three types: positive proposals, incomplete samples, and background samples. We use two types of linear classifiers, a $K + 1$ (K activity classes plus the background) activity classifier predicting the activity category, and K class-specific classifiers determining the completeness of the temporal proposal given that it is a certain

activity. We also employ a location regression scheme similar to practices in [5] to further refine the temporal extent of positive proposals.

For any proposal P_i , with class $K_i \in \{0, \dots, K\}$ and completeness is $C_i \in \{0, 1\}$ if it is not background class, the loss can be formulated as

$$\mathcal{L} = \mathcal{L}_{cls} + \gamma \mathcal{L}_{reg}$$

where

$$\mathcal{L}_{cls} = -\log p(K_i, C_i; P_i) = -\log p(K_i|C_i = 1) - \mathbb{I}_{K_i \geq 1} p(C_i|K_i)$$

The whole pipeline is designed to train in an end-to-end manner. When training, we use sparse sampling strategy similar to TSN [18] such that the memory demand is independent of the proposal length. A strategy resembling online hard negative mining [11] is used to mitigate the imbalance between positive and incomplete samples. When evaluating, the features before STPP is extracted densely with a fixed interval of 6 frames, which make more efficient the calculation of classifier predictions and location regressors afterwards per proposal.

We deploy SSN based on ResNet 101 [6]. Networks pre-trained from Kinetics [9] show consistent improvement over those pre-trained from ImageNet [2].

Table 12. Peformance of using Kinetics [9] as pre-train model. The results are evaluated on ActivityNet 1.2 val set using Inception BN [8].

source	average mAP
ImageNet	26.81%
Kinetics	28.57%

4.2. Unifying Activity and Completeness Classifiers

In additional, we take one step further by unifying the $(K + 1)$ -class activity classifier and a set of K completeness classifiers stated above as one single classifier. Both incomplete and background proposals are treated as negative samples, i.e. the 0-th class , while the positive samples remain the same. We also add a fully-connected layer between the STPP and the final classifier. Although results produced by such a variant are inferior to those of SSN, we find that average mAP can be increased from 29.20% to 29.67% on validation set by doing a fusion of two models.

Table 13. Peformance comparison of using original SSN (model a) and unified SSN (model b). The results are evaluated on ActivityNet 1.3 val set using ResNet-101[6].

method	average mAP
model a (act. + comp.)	29.20%
model b (single classifier)	28.00%
model a+b	29.67%

Table 14. Results on ActivityNet 1.3 [3] test set.

Model	average mAP
BN-Inception, model a	29.34%
ResNet-101, model a+b	31.86%

5. Conclusions

This paper describes our team’s solution to task of untrimmed video classification, trimmed action recognition, and temporal action localization. We propose several 3D spatial-temporal models for action recognition, which is realized by inflating convolutional layers and pooling layers in existing 2D model architectures to 3D version. Models trained from Kinetics [9] are found to be able to effectively transferred to ActivityNet [3], improving classification and detection performance in untrimmed videos.

References

- [1] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017.
- [2] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li. ImageNet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.
- [3] B. G. Fabian Caba Heilbron, Victor Escorcia and J. C. Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970, 2015.
- [4] R. Girshick. Fast R-CNN. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015.
- [5] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [7] H. Idrees, A. R. Zamir, Y.-G. Jiang, A. Gorban, I. Laptev, R. Sukthankar, and M. Shah. The thumos challenge on action recognition for videos in the wild. *Computer Vision and Image Understanding*, 155:1–23, 2017.
- [8] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, pages 448–456, 2015.
- [9] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [10] W. Li, L. Wang, W. Li, E. Agustsson, J. Berent, A. Gupta, R. Sukthankar, and L. Van Gool. Webvision challenge: Visual learning and understanding with web data. *arXiv preprint arXiv:1705.05640*, 2017.
- [11] A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [12] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*, pages 568–576, 2014.
- [13] K. Soomro, A. Roshan Zamir, and M. Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. In *CRCV-TR-12-01*, 2012.
- [14] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, pages 4278–4284, 2017.
- [15] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, pages 2818–2826, 2016.
- [16] D. Tran, L. D. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. In *ICCV*, pages 4489–4497, 2015.
- [17] L. Wang, Y. Xiong, D. Lin, and L. Van Gool. Untrimmednets for weakly supervised action recognition and detection. In *CVPR*, 2017.
- [18] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, pages 20–36, 2016.
- [19] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks for action recognition in videos. *arXiv preprint arXiv:1705.02953*, 2017.
- [20] Y. Xiong, L. Wang, Z. Wang, B. Zhang, H. Song, W. Li, D. Lin, Y. Qiao, L. Van Gool, and X. Tang. Cuhk & ethz & siat submission to activitynet challenge 2016. *arXiv preprint arXiv:1608.00797*, 2016.
- [21] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, D. Lin, and X. Tang. Temporal action detection with structured segment networks. *arXiv preprint arXiv:1704.06228*, 2017.

ActivityNet Challenge 2017: Multi-stream and Multi-task Learning for Video Classification and Action Localization

Gurkirt Singh^{1,2}

Andreas M. Lehrmann²
¹Oxford Brookes University

Fabio Cuzzolin¹

Leonid Sigal²

²Disney Research

gurkirt.singh-2015@brookes.ac.uk, {andreas.lehrmann,lsigal}@disneyresearch.com

Abstract

We present an approach for video classification with fused Inception-V3 network streams and temporal action localization from video classification. For video classification, we make use of multi-task learning comprising simultaneous offset regression and classification in a two-stream architecture. For action localization, we experiment with dynamic programming to generate action proposals from frame-level classification scores. In our experiments, we analyze the behaviour of both streams individually as well as in combination and consider different forms of pre-training.

1. Introduction

Action recognition is a vital precursor for many high-level tasks in computer vision, including real-world scenarios such as video surveillance, video indexing and human-robot-interaction. While early works were focused on traditional tasks like trimmed video classification [12, 18, 19, 8, 10, 16], the rise of convolutional neural networks (CNNs) has also led to tremendous progress in more complex tasks like untrimmed video classification and untrimmed temporal action localization [9, 1], achievements that can be attributed to innovation in terms of hierarchical feature representations [17], the availability of effective regularization techniques [7], and the development of multi-stream architectures [16]. The latter, in particular, have been successfully applied to a wide range of action-related tasks at different levels of granularity [20, 21]. Two-stream features pooled in Fisher vectors and VLAD representations have proven to be very effective in action classification [14] but cannot be directly applied to detection. An orthogonal line of research has considered recurrent neural networks with long short-term memory (LSTM [6]) to obtain temporally smoothed predictions of independent CNN outputs [13].

2. Methodology

At a high level, our approach consists of two parallel convolutional neural networks (CNNs) extracting static (i.e., independent) appearance and optical flow features for each frame. We fuse information from both streams using a convex combination of their respective classification scores to obtain a final result. We will now describe our approach in some more detail.

2.1. Input

Given an input video $V = \{v^{(t)}\}_{t=1}^T$ with frames $v^{(t)}$, we extract independent features for each 3-channel RGB frame at a rate of 24fps by passing normalized color intensities to the RGB stream and a 3-channel optical flow image¹ with (f_x, f_y, m) per pixel to the optical flow stream. Here, (f_x, f_y) encodes the 2D direction and m the magnitude of the flow. Normalization consists of centering and scaling: Images are centered channel-wise by subtracting a fixed mean value of 128 and dividing by 128, resulting in bounded intensities in the interval $[-1, 1]$. Scaling comprises bilinear up- or downscaling such that the shortest image dimension has a size of 340 pixels. We use random crops of size 299×299 for training and the center crop for testing.

2.2. Network Architecture

Both streams use an Inception-V3 network [17] without batch normalization as a base architecture and produce 1024-dimensional feature vectors. We train the overall network using a multi-task loss: (1) Classification: Both streams produce a C -dimensional softmax score vector that is trained using back-propagation with a cross-entropy loss; (2) Regression: In addition to the classification scores, the appearance stream also produces 3-dim. coefficients for each class describing the offset from the boundaries of the

¹We use an approach by Farnebaek [4] to compute the optical flow between two consecutive frames.

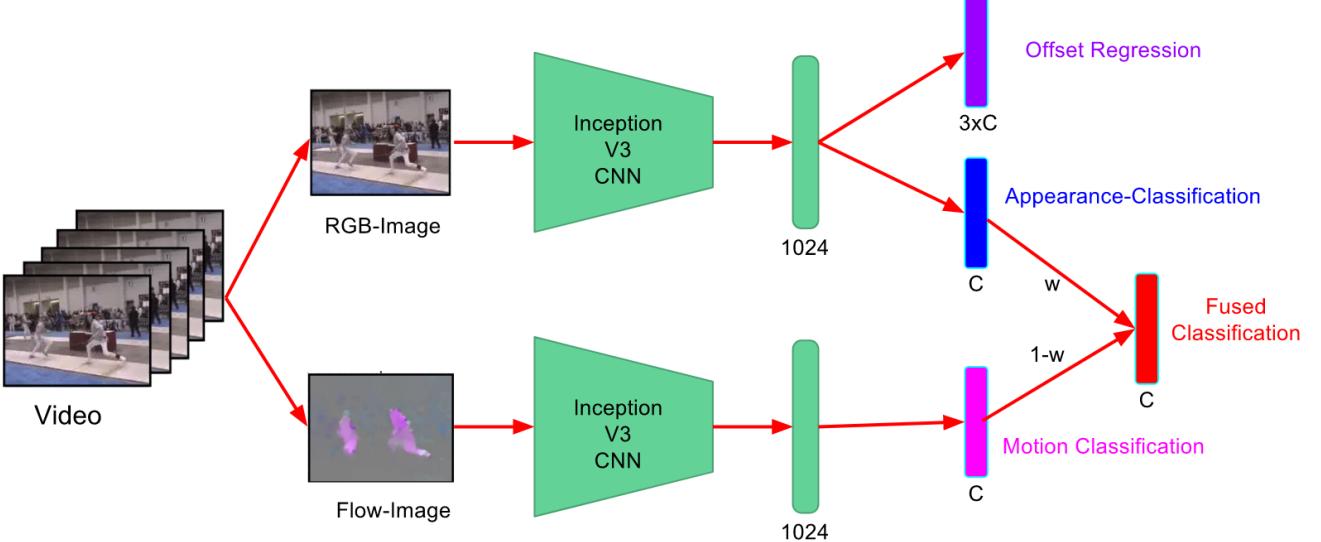


Figure 1: Overview. We show our proposed two-stream model with two parallel Inception-V3 networks [17] processing appearance and motion inputs. While the appearance stream is followed by two sibling layers for offset regression and softmax classification, the motion stream ends in a single softmax classification layer. The final output is obtained as a convex fusion of the classification scores from both streams.

current action as well as its overall duration. This network path is trained using a smooth L_1 loss [5]. Figure 1 summarizes our architecture.

2.3. Fusion

For video classification, we fuse the scores of both streams at the frame-level using a convex combination. The single weight w can be found by cross-validation on the validation set. In practice, we found the values $w = 0.55$ and $1 - w = 0.45$ to be optimal for the appearance and the flow stream, respectively. Finally, we use top- k mean-pooling to obtain a video-level score $s^{(c)}$ for each class c , i.e., $s^{(c)}$ is the mean of the top- k frame-level scores of class c . For a sequence of length T , we found $k = 0.5 \cdot T$ to be a good choice.

2.4. Temporal Action Localization

We generate activity proposals by casting the frame-level scores $S^* = \{s_t^*\}_{t=1}^T$ of the winning class at the video-level (section 2.3) as an optimization problem [3], resulting in disjunct sets of piece-wise constant segments:

Given S^* , we assign a binary label $l_t \in \{0, 1\}$ (where 0 represents the ‘background’ class) to each frame by solving the following optimization problem:

$$\mathcal{C}(L) = \sum_{t=1}^T s_t^* - \lambda \sum_{t=2}^T \psi(l_t, l_{t-1}), \quad (1)$$

where $\mathcal{C}(L)$ is the cost associated with a labeling

$L = \{l_t\}_{t=1}^T$, λ is a scalar weight parameter, and the pairwise potential ψ is defined as

$$\psi(l_t, l_{t-1}) = \begin{cases} 0 & \text{if } l_t = l_{t-1} \\ \alpha & \text{if } l_t \neq l_{t-1} \end{cases}.$$

The cost α for switching from ‘activity’ to ‘background’ is set by cross-validation. Note that this penalizes labelings which are not smooth and thus enforces a piece-wise constant solution. Eq. (1) can be efficiently solved using dynamic programming and is easily extendable to simultaneous detection and classification [3]. To obtain a set of proposals, we treat each positive, contiguous subsequence a in the optimal binary labeling L^* as a proposal whose global score $\mu(a)$ is given by the mean of its corresponding frame-level scores. The final result is given by assigning the label of the highest video-level score $c^* = \arg \max_c s^{(c)}$ (section 2.3) to the top-5 activity proposals and re-weighting each proposal’s score by $\mu(a)' = s^{(c^*)} \cdot \mu(a)$.

3. Experiments

We evaluate our approach on the Kinetics dataset [11] and the ActivityNet 200 dataset [1]. Kinetics consists of $300k$ trimmed videos that are approximately 10 sec. long and labeled with one out of 400 actions. ActivityNet comprises approximately $20k$ videos (split into training/validation/testing at a ratio of $2 : 1 : 1$) falling into 200 activity classes.

We report results for 3 out of 5 ActivityNet challenges:

- (1) Trimmed Action Recognition; (2) Untrimmed Video

Method	Error			
	Top-1	Top-5	Avg	
Val.	Appearance (A)	31.7%	12.7%	22.0%
	Motion (M)	46.4%	23.4%	34.9%
	Fused (A+M)	30.7%	12.3%	21.6%
Test	Fused (A+M)	—	—	22.7%

(a) Trimmed Action Recognition

Val.	Method	Error		
		Top-1	Top-5	Avg
Caba <i>et al.</i> [1]	42.5%			
Appearance* (A*)	25.7%			
Appearance [†] (A [†])	22.7%			
Appearance [‡] (A [‡])	21.8%			
Motion [†] (M [†])	33.1%			
Fused (M [†] +A [†])	20.7%			
Fused (M [†] +A [‡])	19.8%			
Caba <i>et al.</i> [1]	42.2%			
Fused (M [†] +A [‡])	18.8%			

(b) Untrimmed Video Classification

Val.	Method	mAP @ t-IoU threshold δ		
		0.5	0.95	0.5:0.95
With DP	35.5%	1.4%	19.9%	
Without DP	33.2%	4.5%	21.7%	
Test	Without DP	—	—	22.5%

(c) Temporal Action Localization

Figure 2: **ActivityNet Results.** (a) Top-1 error, Top-5 error, and average error for trimmed action recognition on the validation and test set. (b) Top-1 error for untrimmed video classification on the validation and test set. (c) mAP at a temporal intersection-over-union (t-IoU) threshold δ for action localization on the validation and test set. [DP \triangleq Dynamic Programming; (*) Initialized with pre-trained ImageNet weights; ([†]) Initialized with pre-trained Kinetics weights; ([‡]) Initialized with pre-trained Kinetics weights + offset regression.]

Classification; (3) Temporal Action Localization.²³ Our evaluation comprises an ablation study on the public validation set and official results on the withheld test set using the evaluation server [1].

3.1. Trimmed Action Recognition

We train a two-stream model on the Kinetics dataset according to section 2.⁴ Note that this is unlike the stack of images used to train the flow stream in [16]. Both inception networks are initialized with pre-trained ImageNet weights [15]. On the validation set, we also report the performance of the individual streams without fusion to provide further insights. Table 2(a) shows our results. Interestingly, the fused average error on the validation set is 21.6% and thus 0.4% smaller than the error of the appearance stream alone, even though the motion stream alone incurs a rather high average error of 34.9%.

3.2. Untrimmed Video Classification

We initialize the inception network of the flow stream with pre-trained Kinetics weights and experiment with three different training settings for the appearance stream: (1) The appearance stream is initialized with pre-trained ImageNet weights [15]; (2) Similar to [2], the appearance stream is initialized with pre-trained Kinetics weights [11]; (3) In addition to (2) and unlike (1) and (2), we use offset regression as discussed in section 2 and illustrated in Figure 1. We can observe from Table 2(b) that an initial-

ization with Kinetics weights is significantly better than an initialization with ImageNet weights (+3% on the validation set). Using offset regression provides another improvement of 0.9%. Fusion with the motion stream leads to yet another performance gain of 2%, reaching a final top-1 error of 19.8%.

3.3. Temporal Action Localization

We explore two simple techniques of achieving activity localization based on video-level scores: (1) The activity proposals are generated using dynamic programming (DP), as described in section 2.4; (2) Each video produces a single activity proposal that lasts from the start of the video to its end. Scoring of the proposals follows our description in section 2.4. Although DP helps to break the video sequence into smooth activity proposals and achieves a good performance at a t-IoU threshold of $\delta = 0.5$, we found that it typically hurts the performance for higher thresholds, as shown in Table 2(c). Surprisingly, a naive activity proposal from start to end is therefore better than DP for higher and/or integrated thresholds, which is shown in the second row of Table 2(c).

4. Conclusion and Future Work

We’ve presented an ablation study on activity recognition that highlights the importance of adequate pre-training, multi-task learning, and proper fusion. Furthermore, we’ve discussed two simple approaches to achieve action localization and have shown that video-level classification is beneficial for activity localization.

²(1) is evaluated on Kinetics, (2) and (3) are evaluated on ActivityNet.

³Official results for (2) and (3) were not available at the time of submission and will be included in the final version of this manuscript.

⁴Different from the base architecture described in section 2, we do not use offset regression for trimmed action recognition.

References

- [1] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015. [1](#), [2](#), [3](#)
- [2] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. *arXiv preprint arXiv:1705.07750*, 2017. [3](#)
- [3] G. Evangelidis, G. Singh, and R. Horaud. Continuous gesture recognition from articulated poses. In *ECCV Workshops*, 2014. [2](#)
- [4] G. Farnebäck. Two-frame motion estimation based on polynomial expansion. *Image analysis*, pages 363–370, 2003. [1](#)
- [5] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1440–1448, 2015. [2](#)
- [6] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. [1](#)
- [7] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015. [1](#)
- [8] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(1):221–231, Jan 2013. [1](#)
- [9] Y. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. Thumos challenge: Action recognition with a large number of classes. <http://crcv.ucf.edu/THUMOS14>, 2014. [1](#)
- [10] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *IEEE Int. Conf. on Computer Vision and Pattern Recognition*, 2014. [1](#)
- [11] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. [2](#), [3](#)
- [12] I. Laptev, M. Marszałek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *IEEE Int. Conf. on Computer Vision and Pattern Recognition*, 2008. [1](#)
- [13] S. Ma, L. Sigal, and S. Sclaroff. Learning activity progression in lstms for activity detection and early detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1942–1950, 2016. [1](#)
- [14] A. Miech, I. Laptev, and J. Sivic. Learnable pooling with context gating for video classification. *arXiv preprint arXiv:1706.06905*, 2017. [1](#)
- [15] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. [3](#)
- [16] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in Neural Information Processing Systems 27*, pages 568–576. Curran Associates, Inc., 2014. [1](#), [3](#)
- [17] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016. [1](#), [2](#)
- [18] H. Wang, A. Kläser, C. Schmid, and C. Liu. Action Recognition by Dense Trajectories. In *IEEE Int. Conf. on Computer Vision and Pattern Recognition*, 2011. [1](#)
- [19] H. Wang and C. Schmid. Action Recognition with Improved Trajectories. In *Proc. Int. Conf. Computer Vision*, pages 3551–3558, 2013. [1](#)
- [20] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *European Conference on Computer Vision*, pages 20–36. Springer, 2016. [1](#)
- [21] Z. Yuan, J. C. Stroud, T. Lu, and J. Deng. Temporal action localization by structured maximal sums. *arXiv preprint arXiv:1704.04671*, 2017. [1](#)

LNMIIT Submission to the ActivityNet Challenge 2017

Arpan Gupta

The LNM Institute of Information Technology

Jaipur, India

arpan@lnmiit.ac.in

Abstract

In this paper, a brief description is provided of the method that was used for the tasks of untrimmed video classification, temporal action proposals and temporal action localization. The C3D features were used to train a fully connected network for the 200 action categories. The trained model was used to obtain probabilities over all the categories for a sequence of validation features of a video. The action category with the highest weightage of votes was the predicted action class. From the sequence of probability values over all the classes, the temporal proposals and action localization results were obtained after tuning a few thresholding parameters on the validation set. The above method gave a **Top-1 error of 0.36839** in untrimmed classification, **AUC of 39.8498** for temporal action proposals, and **Avg. mAP of 0.13182** for temporal action localization task.

1. Introduction

The ActivityNet Challenge 2017 had five tasks, namely, Untrimmed Video Classification, Trimmed Action Recognition, Temporal Action Proposals, Temporal Action Localization and Dense-Captioning Events in Videos (hereafter, denoted by T1, T2,..., T5). The tasks T1, T3, T4 and T5 were based on the ActivityNet dataset [3] while T2 used Kinetics dataset [7]. These datasets are the state-of-the-art benchmark datasets for the video activity recognition tasks. As such, they have a large number of annotated videos for 200 classes in ActivityNet and 400 classes in Kinetics.

This work describes the methodology followed for my submission to tasks T1, T3 and T4, which involve only the ActivityNet dataset. Task T1 required to predict the activity categories for the untrimmed videos, where the video may belong to more than one category. For task T3, action proposals had to be generated, on the lines of [4] such that the area-under-the-curve(AUC) for Avg. Recall Vs Avg. Number of Proposals was maximized. Task T4 was the action localization/detection task, where one had to predict in which

part of a video a specific action was occurring.

In this report, section 2 details the training and prediction steps, section 3 covers the results and related discussion, while section 4 gives concluding remarks for the work carried out.

2. Methodology

The C3D features [11] were used to train a fully connected network in Caffe [5]. Code available¹. The C3D features were provided by the organizers. They were the outputs of the FC-7 layers (4096 dimensional vectors) for every eighth frame using a pre-trained model trained on sports-1M [6] clips of 16 frame temporal resolution. Further, PCA was applied to reduce their dimensionality from 4096 to 500.

Global video features, such as the ImageNet Shuffle features [9] and the MBH features [12], can be used for T1, involving global prediction, but cannot be reliably used for T3, T4, and T5 where a frame-level classification is needed. Therefore, the C3D features were chosen out of the provided features.

2.1. Training

Steps for sampling and training:

1. Randomly select 4k unique positions (frame numbers) from actions of each category. For 200 classes of ActivityNet we get 800k training samples(1k per class for validation set). Note that the samples are shuffled.
2. Obtain C3D feature vector (each of size 500) for each sample by dividing by 8 (since, the C3D features are obtained for every eighth frame). Take care of boundary cases by subtracting by 1.
3. Read the feature vector from h5py database and saved into LMDB with the corresponding label assigned. The creation of LMDB is done parallelly, by incrementally saving in batches of 800 vectors at a time. Paral-

¹Github link: www.github.com/arpane4c5/ActivityNet

Layer	#Outputs	#Parameters
FC-1	1024	500×1024
ReLU	1024	0
Dropout	Ratio=0.5	0
FC-2	1024	1024×1024
ReLU	1024	0
Dropout	Ratio=0.5	0
FC-3	200	1024×200

Table 1. Fully-Connected(FC) network. Total Params=1765376.

lization ensured speedup and in case of failure, one may continue from the last written batch.

4. Prepare fully connected network by defining the details for all the layers. The details of the network are given in Table 1
5. The network is trained using Stochastic Gradient Descent for 400k iterations. The stepsize was 100k with base learning rate(α) = 0.01, momentum(μ) = 0.9, and drop factor of 10 (γ = 0.1). The weights of the FC layers were initialized using Xavier initialization [2].
6. The progress of loss for training and validation sets is shown in Fig. 1.

The model was trained on an NVIDIA GeForce GTX 750 card which has a memory of 1GB. This GPU does not have enough capacity to support training/finetuning of larger models, such as C3D models, which is left as a future work.

The model achieved an accuracy of **47.39%** on the validation set samples.

2.2. Predictions

The trained model was used to get prediction probabilities for all the C3D features of the validation set videos. For a single video C3D features, the below steps were followed to come up with the predictions for T1, T3 and T4:

1. For a video $v_i \in V_{val}$, the prediction probabilities for a sequence of C3D features are $p_t^i \in \mathbb{R}^{200}$, 200 is the number of categories in ActivityNet.
2. The p_t^i s are summed up for the video v_i and the 3 highest value categories are predicted as the result, along with the corresponding summed up probability values as the score.
3. Model parameters and number of classes to be predicted was selected using the validation set.
4. The Top-1 error in task T1 was **0.36839**.

Task T3 and T4 needed the temporal extent in the video and a model able to make frame-wise predictions. Therefore, a list of category predictions is obtained corresponding to each C3D feature vector. This list of temporal predictions is of the same length as the total number of C3D feature vectors for a video. We might have top- n category predictions, therefore, for each of the top- n categories, we analyze the temporal segment of that activity in the list. A longer temporal extent means that the model is more confident of occurrence of that activity within the time interval. e.g. if list is [9, 9, 9, 8, 6, 6, 9, 1, 1, 1, 1, 1], the top-3 predicted categories are 9, 1, and 6 (Note that summed up probabilities of category 9 is greater than that of category 1). We first choose 9 and analyze the inter-segment distances (int_seg_dist). Here the int_seg_dist is 4 (between first segment of three 9's and second segment of one 9). Similarly, we can find the int_seg_dist values for each pair of neighbours and if it is below a threshold then merge the neighbouring segments.

For T3, the two parameters Top- n and int_seg_dist were chosen as 12 and 150, respectively, by grid search using the validation set. For T4, int_seg_dist was 60 and n was 1. As the features are for every eighth frame of the video, a more accurate temporal prediction is not possible using C3D features. Therefore, a random value (sampled from a standard gaussian) is added to the starting time and ending time of the proposals.

3. Results and Discussion

The fully-connected network trained on C3D features, described in section 2.1 had a validation accuracy of 47.39%. This accuracy is quite low, and further analysis is needed in this regard. Even increasing the sample size from 4k per class to 10k per class was tried but it did not push the accuracy, significantly. Also, Adam [8] solver method was tried but that too gave similar results.

The **Top-1 error** of **0.36839** was obtained for T1, **AUC** of **39.8498** for T3, and **Avg. mAP** of **0.13182** for T4 on the test set predictions.

The described frame-wise prediction model may give better results if the training is performed on more distinguishable video features. The training of an Alexnet-like model on sample images (resized to 120×160) taken from random positions (as explained in section 2.1) was also tried, such that we get predictions for each frame of a video, but the accuracy did not increase beyond 21% on validation samples. Similarly, optical flow [1] visualizations were used for training an Alexnet-like model, but nearly half of the samples gave blank frames (all zero pixel values). Even increasing the distance between two frames for obtaining the optical flow does not help. Though, the number of samples with all zeroes reduced but still about one-quarter of the samples were zeroes.

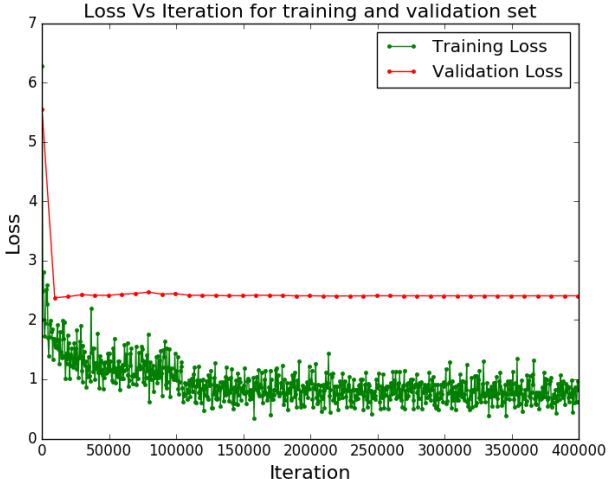


Figure 1. Losses Vs Iterations. Batch Size=64. Epochs=32.

Another idea that was tried was to train fully-connected network on 2D HOG [10] features (9576 vector size for an image of 120×160 pixels). In this case also, the training loss decreased but the testing loss did not and the accuracy even after 120k iterations was constant at 11% on the validation samples.

The global features such as ImageNet shuffle and MBH features had one feature vector for each video, thus, making them suitable for T1. Separate random forest classifiers were trained on only the ImageNet shuffle features for each activity, using all positive samples and randomly selecting same number of negative samples. The submission result Top-1 error for T1 was 0.45585. The number of trees for random forest classifiers were selected by using the validation set and were set to 200 for each category classifier. Using MBH features with random forests, resulted in a worse result than using ImageNet shuffle features. A detailed analysis in this regard is still sought.

4. Conclusion

This submission, though, was not among the top submissions, but it shows that the frame-level features such as C3D features are relatively better features than the global features such as ImageNet Shuffle and MBH features, for the tasks of action classification and temporal action localization. Being frame level features, they can be used for more than a single task by manipulating the frame level prediction probabilities.

References

- [1] G. Farnebäck. Two-frame Motion Estimation Based on Polynomial Expansion. In *Proceedings of the 13th Scandinavian Conference on Image Analysis, SCIA'03*, pages 363–370, Berlin, Heidelberg, 2003. Springer-Verlag.
- [2] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 9:249–256, 2010.
- [3] F. C. Heilbron, V. Escorcia, B. Ghanem, and J. C. Niebles. ActivityNet: A large-scale video benchmark for human activity understanding. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 07-12-June:961–970, 2015.
- [4] F. C. Heilbron and J. C. Niebles. Fast Temporal Activity Proposals for Efficient Detection of Human Actions in Untrimmed Videos. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1914–1923, 2016.
- [5] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional Architecture for Fast Feature Embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [6] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1725–1732, 2014.
- [7] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman. The Kinetics Human Action Video Dataset. 2017.
- [8] D. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]*, pages 1–15, 2014.
- [9] P. Mettes, D. C. Koelma, and C. G. M. Snoek. The ImageNet Shuffle: Reorganized Pre-training for Video Event Detection. 2016.
- [10] B. T. N. Dalal. Histograms of Oriented Gradients for Human Detection. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005.
- [11] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning Spatiotemporal Features with 3D Convolutional Networks. *Iccv*, 2015.
- [12] H. Wang and C. Schmid. Action Recognition with Improved Trajectories. *Computer Vision (ICCV), 2013 IEEE International Conference on*, pages 3551–3558, 2013.

Joint Classification and Detection Using LSTMs

UCSB & IBM Research Submission to ActivityNet Challenge 2017

Oytun Ulutan^{1,2}, Swati Rallapalli², Mudhakar Srivatsa² and B.S. Manjunath¹

¹University of California, Santa Barbara

²IBM Research

Abstract

This paper describes the UCSB & IBM Research submission to ActivityNet 2017 [1] Challenge for Task1: Untrimmed Video Classification (Classification) and Task4: Temporal Action Localization (Detection) tasks. We train a Long-Short Term Memory (LSTM) model on C3D [6] features extracted from 16-frame clips. The model is trained by optimizing on multiple loss functions for both tasks jointly. Results indicate that optimizing on this multi-task loss improves performance. In addition to deep features, to improve the classification accuracy, we leverage the video-level hand crafted Imagenet-Shuffle [4] and Motion Boundary Histogram [7] features and train one-versus-rest linear SVMs for each activity class and each feature. Final classification results are obtained by averaging probabilities from LSTM and both SVMs. The resulting classification label is assigned to the action region obtained from the detection result of LSTM model.

1. Introduction

Action classification and detection are important and challenging tasks in computer vision. To temporally detect the action, methods are needed to model the temporal evolutions of the actions over time in addition to spatial appearance. This requires models to learn actions that span from a couple of frames to thousands of frames. Learning such models require large sets of annotated videos. ActivityNet 2017 Challenge [1] provides 648 hours of annotated videos of 200 action classes. Our proposed model is tested on Task1: Untrimmed Video Classification (Classification) and Task4: Temporal Action Localization (Detection) tasks.

2. Methodology

Our proposed method uses both global-features and frame-level features to obtain the final Classification and

Detection results. Linear SVMs are trained for each of the global-features. A Long-Short Term Memory (LSTM) model is trained jointly for Classification and Detection.

2.1. Global-Features for Classification

Imagenet-Shuffle [4] and Motion Boundary Histograms [7] features provided by the ActivityNet¹ are used.

Imagenet-Shuffle features are global-features that are generated by GoogleNet model [5] trained on Imagenet. Convolutional features are extracted from each frame and mean pooled across the frames to obtain a single global descriptor for each video.

Motion Boundary Histograms are extracted using Improved Trajectories [7] method. These histogram descriptors encode the motion between pixels and generate a global-feature vector for each video.

For each global-feature one-versus-rest Linear SVM classifiers are trained for the task and the scores are scaled into probabilities using a softmax function.

2.2. Frame-Level Joint Classification and Detection

Using frame-level C3D features, LSTM model is trained while optimizing jointly on multiple loss functions. We propose and experiment on different regularization methods for videos and report the results.

C3D features are frame-level features that are extracted using the C3D model [6] on Sports 1M dataset [3]. The features are not fine-tuned for the ActivityNet dataset. C3D model extracts 3D Convolutional features on 16-frame inputs. The features are extracted for every 8 frames, i.e., with 50% overlap. Furthermore, the dimensionality of the features are reduced from 4096 to 500 dimensions using PCA. These features are also provided by ActivityNet organizers.

LSTMs have the ability to model the evolution of actions over time by leveraging memory cells. In our model we used the LSTM cell implementation from [2]. Figure

¹<http://activity-net.org/challenges/2017/download.html>

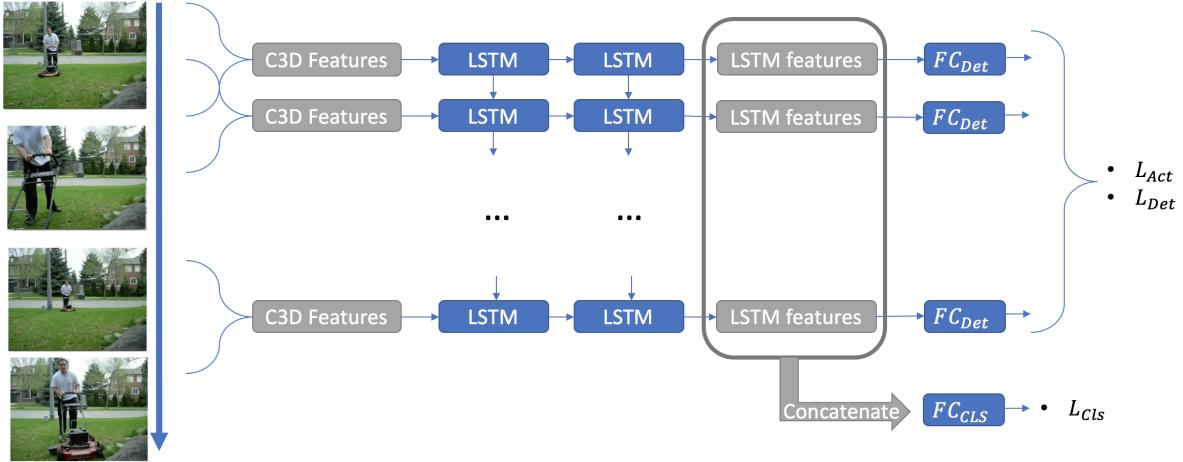


Figure 1. Frame-Level LSTM model for joint Detection and Classification. Model is jointly trained for both Classification and Detection tasks using a Classification loss L_{cls} , an Actionness loss L_{act} and a Detection loss L_{det} .

1 shows our frame-level LSTM model architecture. LSTM model is trained for both detection and classification tasks. This is achieved by jointly optimizing on three different loss functions, namely, a Classification loss L_{cls} , an Actionness loss L_{act} and a Detection loss L_{det} .

2.2.1 Detection (L_{det}) and Actionness (L_{act}) Losses

Outputs of the LSTM cells at each timestep individually get processed by a Fully Connected layer to obtain a sequence of 201 dimensional class (200 action classes and a background class) probabilities. These are used to calculate frame-level detection cross entropy loss L_{det} and frame-level actionness cross entropy loss L_{act} .

When calculating the frame-level losses, we impose different weights for background and action classes. This allows us to penalize Action false negatives more than Background false negatives and also helps with class imbalance. We choose these weights to be 0.5 for Background and 1 for Action classes.

Additionally, to achieve better accuracy in detecting start and end points of the actions, we impose a transition weight for the sequences where actions start or end. This allows us to penalize the system more on action transitions and allows the model to learn the transitions more robustly. We choose these weights to be 0.5 for no transition sequences and 1 for transition sequences.

2.2.2 Classification loss (L_{cls})

For the classification task, LSTM features over a defined number of timesteps are concatenated together and a Fully Connected layer is used to obtain a single class label for

Method	mAP(0.50:0.05:0.95)	Hit@1 Acc
Imagenet-Shuffle SVM	-	63%
MBH-Shuffle SVM	-	59%
SVMs (Scores Averaged)	-	64%
$L_{act} + \text{SVMs}$	0.12	64%
$L_{act} + L_{det} + L_{cls}$	0.14	63%
$L_{act} + L_{det} + L_{cls} + \text{SVMs}$	0.15	67%

Table 1. Results on Validation set of ActivityNet Challenge. Detection performance is measured by calculating the average mAP over temporal Intersection over Union (tIoU) thresholds (0.50:0.05:0.95). For Classification Hit@1 accuracy is calculated for each video. These are the metrics used for the leaderboard on Testing set.

that time sequence. We define the correct label for this sequence as the most frequent ground truth label in that interval. These labels, compared to using video-level labels, allows us to have background instances in the time intervals with no action happening. These labels are used to calculate the video-sequence level classification loss L_{cls} .

3. Experiments

We test the performance of each module by training the models with proposed loss functions, one-by-one, to evaluate the benefit of each. For the performance metric, we follow the ActivityNet 2017 Challenge guidelines. For the temporal detection task, mean Average Precision (mAP) values are averaged over temporal Intersection over Union(tIoU) thresholds in between 0.50 and 0.95 (inclusive) with a step size of 0.05. For the classification task, we report the accuracy of the highest scored class label for each video. Table 1 demonstrates the results from our ex-

periments.

References

- [1] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
- [2] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [3] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- [4] P. Mettes, D. Koelma, and C. G. M. Snoek. The imagenet shuffle: Reorganized pre-training for video event detection. In *Proceedings of the ACM International Conference on Multimedia Retrieval*, New York, USA, 2016.
- [5] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [6] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015.
- [7] H. Wang and C. Schmid. Action recognition with improved trajectories. In *IEEE International Conference on Computer Vision*, Sydney, Australia, 2013.

NTHU CVLab at ActivityNet Challenge 2017 - Task 2: Trimmed Action Recognition

Yi Huang*, Hsing Yu Chen*, Wei Jing Wong, and Shang-Hong Lai

Department of Computer Science
National Tsing Hua University

{jeffreyhuang0823, andy19933, wjwilliamwong}@gmail.com, lai@cs.nthu.edu.tw

Abstract

We implement a multi-stream framework to utilize the rich multi-modal information in videos for human action recognition. Specifically, we first train three convolutional neural networks whose inputs are RGB images, stacked optical flow and human poses to model spatial, short-term motion and pose information respectively in each video.

1. Introduction

Encouraged by the success of two-stream Convolutional Neural Network (CNN) [1] on action recognition in recent years, many researchers have developed some outstanding works based on this structure such as TSN [2] and Temporal-Inception [3]. Moreover, a few recent studies attempted to introduce additional stream to model extra clues in videos [4]. In this work, we propose to add pose information into our network to evaluate whether the pose stream can further enhance the performance of multi-stream CNNs.

2. Method

We implement a three-stream model to utilize spatial, short-term motion and pose information in videos to recognize actions in them. In this section, we describe the implementation details of the three streams. The overview of our pipeline is illustrated as in Figure 1.

The spatial stream directly models spatial or pixel information from visual frames. With the outstanding performance proven on ImageNet, we expect the frame-level features to also yield high performance in action recognition. On the other hand, the motion stream models short-term motion information in videos through stacked optical flow images. The motion stream complements the spatial stream so that the videos are not only classified based on static

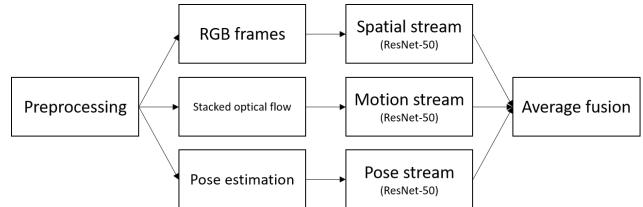


Figure 1. The overview of our pipeline.

frames but also movements of humans and objects between frames. Furthermore, the pose stream omits the contextual and spatial information of the video scene while focusing on the posture of human bodies. It derives the actions the people in the video are performing based on their poses, which may consist of joints and rigid body parts.

2.1. Data Preprocessing

2.1.1 Video Frames Extraction

Since all of our input to spatial, motion, and pose stream are based on video frames, we extract frames from Kinetics dataset using two sampling strategies.

The first sampling strategy, we extract frames with sampling rate of $\frac{1}{10}$ and resize them to 256×256 . For example, if a video has 300 frames, f_0, f_1, \dots, f_{299} , using sampling rate $\frac{1}{10}$ we only save $f_0, f_{10}, f_{20}, \dots, f_{280}, f_{290}$. This sampling strategy greatly reduces the number of training data by a factor of 10. Therefore, it provides a trade-off between preprocessing time and the amount of information we are able to utilize, which may then affect the final performance.

As for the second sampling strategy, we extract every frame in the videos and resize them to 256×256 . This version is mainly for the optical flow computation phase.

2.1.2 Optical Flow Computation

We use Flow-net 2.0 docker version [5] to compute optical flow images based on the video frames we extracted. Since the number of frames is very

*These authors contribute equally to this work.

huge, we only compute optical flow once every 10 frames. For example, if a video has 300 frames, f_0, f_1, \dots, f_{299} , we compute optical flow on frame pairs $(f_0, f_1), (f_{10}, f_{11}), (f_{20}, f_{21}), \dots, (f_{280}, f_{281}), (f_{290}, f_{291})$ only. Again, this is a trade-off between preprocessing and training time and performance.

2.1.3 Human Pose Estimation

In order to address the appearance of multiple people in the videos, we use the multi-person pose estimation method [6] to extract human poses. The human pose is represented by 15 joints of human body. The pose estimation process outputs 15 confidence maps corresponding to the 15 joints, indicating the confidences of joints detected at certain locations. These confidence maps are then aggregated as a 15-channel input to the pose stream of the proposed framework.

2.2. CNN Architecture

Our three-stream model is constructed by three individual spatial, motion and pose stream. Each stream models different type of information in videos respectively and independently. Due to the large size of Kinetics dataset, and some hardware resource limitation, we choose ResNet-50 as the backbone of our networks which has a considerably good balance between accuracy and training speed. The three ResNet-50 CNNs are initialized with weights pre-trained on ImageNet and fine-tuned with Kinetics dataset. In order to utilize the ImageNet RGB models to initialize the motion and pose model, we use cross-modality pre-training method [2] that averages the weights across RGB channels and replicates this average by the channel number of the target network.

2.3. Fusion

We choose the averaging fusion approach as it is a simple and fast approach. For a given test video, we assign the prediction scores from the spatial, motion and pose stream with pre-defined weightings and average them together.

3. Experiments

We report the performance of each stream and the fused result on the validation set of Kinetics dataset in Table 1, where each of the three streams is trained individually and fused to obtain the results. We can see that the fused result achieves the best performance compared to any individual stream. Note that due to some downloading problems, there are about $\frac{1}{10}$ to $\frac{1}{12}$ missing videos in our downloaded Kinetics training and validation sets; thus all the reported figures in Table 1 do not represent the actual results on the full dataset released by the organizer.

We also compare different fusion combinations of the three streams to see how much improvement each of them can bring. The results are reported in Table 2. We can see that fusing all the three streams does not end with the best average error rate. Instead, fusing the spatial and motion streams does.

Stream	Top-1 Err.	Top-5 Err.	Avg. Err.
Spatial	0.3157	0.1485	0.2321
Motion	0.7173	0.4994	0.6084
Pose	0.7697	0.5658	0.6678
Fused(all)	0.2956	0.1403	0.2180

Table 1. Performance on the validation set of Kinetics

Combination	Top-1 Err.	Top-5 Err.	Avg. Err.
S + M	0.2950	0.1406	0.2178
S + P	0.3140	0.1469	0.2305
M + P	0.6851	0.4483	0.5667
Fused(all)	0.2956	0.1403	0.2180

Table 2. Ablation analysis of different fusion combinations

4. Conclusions and Future Work

In a nutshell, we proposed a three-stream CNN exploiting the spatial, motion and human pose information of the trimmed videos for human action recognition. The fusion of these information has been proven to outperform individual streams. In order to achieve state-of-the-art performance, the architecture of the CNN used in each of the three streams can be further improved to address different properties in the data.

References

- [1] K. Simonyan and A. Zisserman, “Two-stream convolutional networks for action recognition in videos,” in *NIPS*, 2014.
- [2] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. V. Gool, “Temporal segment networks: Towards good practices for deep action recognition,” in *ECCV*, 2016.
- [3] Z. K. G. A. Chih-Yao Ma, Min-Hung Chen, “Ts-lstm and temporal-inception: Exploiting spatiotemporal dynamics for activity recognition,” in *CVPR*, 2017.
- [4] Z. Wu, Y.-g. Jiang, X. Wang, H. Ye, and X. Xue, “Multi-stream multi-class fusion of deep networks for video classification,” 2016.
- [5] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, “Flownet 2.0: Evolution of optical flow estimation with deep networks,” in *CVPR*, 2017.
- [6] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, “Realtime multi-person 2d pose estimation using part affinity fields,” in *CVPR*, 2017.

DevNet2.0: Revisiting the Effectiveness of Off-shelf modeling approaches for Large-scale Video Classification

Yunlong Bian, Chuang Gan, Xiao Liu, Fu Li, Xiang Long
Yandong Li, Heng Qi, Jie Zhou, Shilei Wen, Yuanqing Lin
Baidu IDL & Tsinghua University

Abstract

This paper describes our solution for the video recognition task of ActivityNet Kinetic challenge that ranked the 1st place. Most of existing state-of-the-art video recognition approaches are in favor of an end-to-end pipeline. One special case is the framework of DevNet [2]. The merit of DevNet is that they first use the video data to learn a network (i.e. fine-tuning or training from scratch). Then they feed the features from the learned network into the off-shelf machine models to conduct video classification instead of directly using the end-to-end classification score (e.g. softmax score). However, the effectiveness of this line work has long-term been ignored and underestimated. In this submission, we extensively follow this strategy. Besides, we also propose four novel temporal modeling approaches using the learned features: Multi-group Shifting Attention Network, Temporal Xception Network, Multi-stream sequence Model and Fast-Forward Sequence Model. Experiment results on the challenging Kinetic dataset demonstrate that our proposed temporal modeling approaches can significantly improve existing approaches in the large-scale video recognition tasks. To be noted, our best single all attention network achieves 78.7% in top 1 and 93.9% on the validation set.

1. Introduction

Video understanding is among one of the most fundamental research problems in computer vision and machine learning. The ubiquitous video acquisition devices (e.g., smart phones, surveillance cameras, etc.) have created videos far surpassing what we can watch. It has therefore been a pressing need to develop automatic video understanding and analysis algorithms for various applications.

To recognize actions and events in videos, recent approaches based on deep convolutional neural networks (CNNs) [5, 9, 2, 13] and/or recurrent networks [4, 11, 1] have achieved state-of-the-art results. However, due to the lack of public available datasets, existing video recognition

approaches are restricted to understand small-scale data, while large-scale video understanding remains an under-addressed problem. To remedy this issue, Google DeepMind releases a new large-scale video dataset, named as Kinetic dataset which contains 300K video clips of 400 human action class.

To address this challenge, our solution follow the strategy of DevNet framework [2]. Particularly, we first learn the RGB, Flow and Audio neutral network model using the video. Then we extract the multi modality feature and fed them into different off-shelf temporal models. We also propose four novel temporal modeling approaches, namely Multi-group Shifting Attention Network, Temporal Xception Network, Multi-stream sequence Model and Fast-Forward Sequence Model. Experiment results verity the effectiveness of the four models over the traditional temporal modeling approaches. We also find that these four temporal modeling approaches are complementary with each others and lead to the state-of-the-arts performances after ensemble.

The remaining sections are organized as follows. Section 2 presents the basic multimodal feature extraction. Section 3 describe our proposed off-shelf temporal modeling approach. Section 4 reports empirical results, followed by discussion and conclusion in Section 5.

2. Multimodal Feature Extraction

Videos are naturally multimodal because a video can be decomposed into visual and acoustic components, and the visual component can be further divided into spatial and temporal parts. We extracted multimodal features to best represent videos accordingly.

2.1. Visual Feature

As in [9], we used RGB images for spatial feature extraction and used stacked optical flow fields for temporal feature extraction. We tried different ConvNet architectures and found Inception-ResNet-v2 [12] outperforms others in both spatial and temporal components. The RGB model is

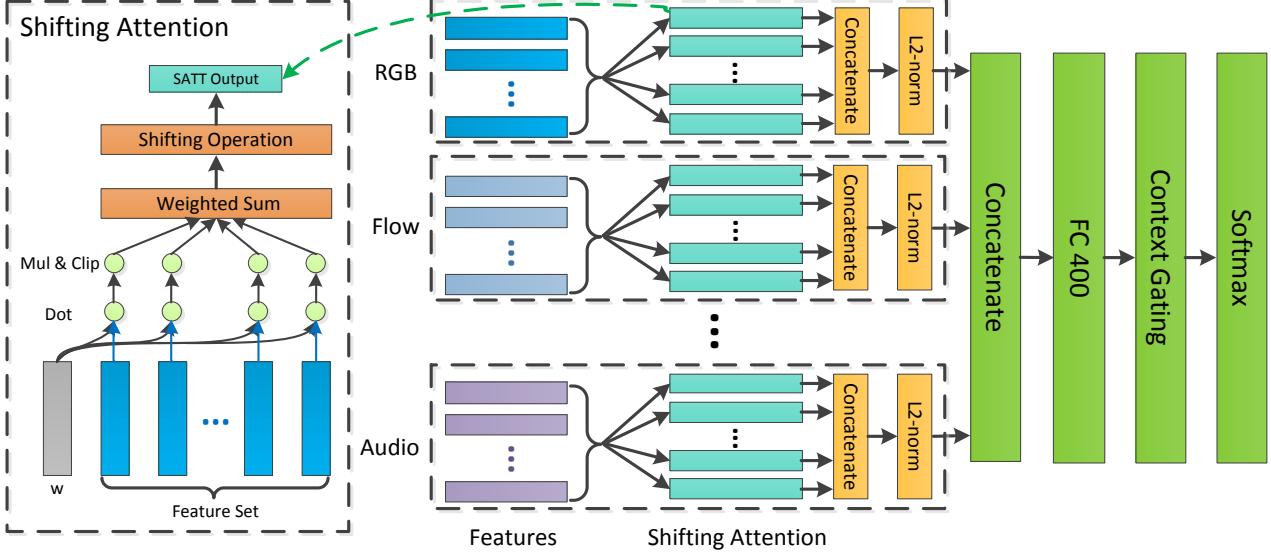


Figure 1. Multi-group Shifting Attention Networks.

initialized with pre-trained model from ImageNet and fine-tuned in the Kinetics dataset, while the flow model is initialized from the fine-tuned RGB model. Inspired by [15], the temporal segment network framework is used and three frames are sampled from each trimmed video for video-level training. During testing, we can densely extract features for each frames in the video.

2.2. Acoustic Feature

We use ConvNet-based audio classification system [3] to extract acoustic feature. The audio is divided into 960ms frames, and the frames are processed with Fourier transformation, histogram integration and logarithm transformation. The resulting frame can be seen as a 96×64 image that form the input of a VGG16 [10] image classification model. Similar with the visual feature, we trained the acoustic feature in the temporal segment network framework.

3. Off-shelf Temporal Modeling Approaches

We present our proposed shifting attention network, temporal Xception network in this section. We refer [6] for the implementation details of multi-stream sequence model and fast-forward sequence model.

3.1. Shifting Attention Network

Recently, attention model shows great potentials in sequence modeling, e.g., many pure attention architectures [14, 7] are proposed and achieve promising results in natural language processing problems. In order to explore the capabilities of attention models in action recognition, an shifting attention network architecture is proposed, which is efficient, elegant and solely based on attention.

3.1.1 Shifting Attention

An attention function can be considered as mapping a set of input features to one output, where the input and output are both matrixes that concatenate feature vectors. The output of shifting attention $SATT(X)$ is calculated through a shifting operation based on weighted sum of the features:

$$SATT(X) = \frac{\lambda X \cdot a + b}{\|\lambda X \cdot a + b\|_2}, \quad (1)$$

where λ is a weight vector calculated as

$$\lambda = \text{softmax}(\text{clip}_{-\beta}^{\beta}(\alpha \cdot w X^T)), \quad (2)$$

w is learnable vector, a and b are learnable scalars, and α and β are hyper-parameters to control the sharpness of the distribution. The shifting operation actually shift the weighted sum and at the same time ensuring the invariance of the scale. The shift operation efficiently introduces diversity for each attention, laying the foundation for Multi-SATT.

3.1.2 Multi-group Shifting Attention Networks

In order to collect multimodal information from videos, we extract a variety of different features, such as appearance (RGB), motion (flow) and audio. Although attention focus on some specific features and effectively filters irrelevant noise, it's unrealistic to include multimodal feature sets within one attention. Here, we propose a Multi-group Shifting Attention Networks for training multiple groups of attentions simultaneously. The architecture of the proposed Multi-SATT is illustrated in Figure 1.

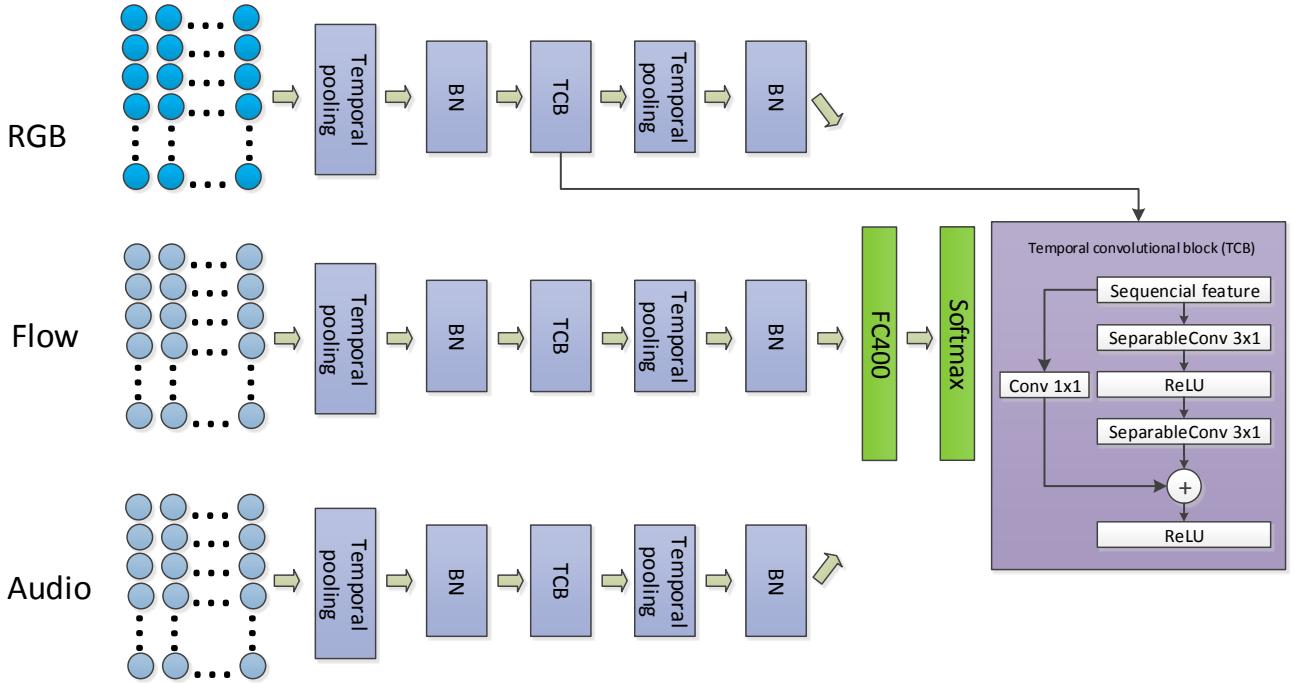


Figure 2. Temporal Xception Network.

First, we extracted multiple feature sets from the video. For each feature set X_i , we applied N_i different shifting attentions, which we call one attention group, and concatenated the outputs. Next, outputs of different attention groups were normalized separately and concatenated as an overall representation vector of the video. Finally the representation vector was used for classification through a fully-connected layer. Context Gating [8] was used before softmax layer.

3.2. Temporal Xception Network

Depthwise separable convolution architecture has shown its power in image classification by reducing the number of parameters and increasing classification accuracy simultaneously. Recently, convolutional sequence-to-sequence networks have been applied to machine translation tasks. In this competition, we introduce temporal Xception network for action recognition, which is in the depthwise separable convolution families and achieves promising performance.

The proposed temporal Xception network architecture is shown in Figure 2. Zero-valued multimodal features were padded to make fixed length data for each stream. We applied adaptive temporal max pooling to obtain n segments for each video. We then propagated the batch data into a Temporal Convolutional block, which is a stack of 2 Separable convolutional layers connected by batch norm and activation with a shortcut connection. Finally, the outputs of three stream features are concatenated and fed into the fully-

connected layer for classification.

4. Experiment Results

The Kinetic dataset contains 246535 training images, 19907 validation images and 38685 testing images. Each image is in one of 400 categories. Table 1 summarizes our results on the Kinetics validation dataset.

5. Conclusions

In this work, we have proposed four temporal modeling approaches to address the challenging large-scale video recognition task. Experiment results verify that our approaches achieve significantly better results than the traditional temporal pooling approaches. The ensemble of our individual models has been shown to improve the performance further, enabling our method to rank first worldwide in the challenge competition. All the code and models will be released soon.

References

- [1] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
- [2] C. Gan, N. Wang, Y. Yang, D.-Y. Yeung, and A. G. Hauptmann. Devnet: A deep event network for multimedia event detection and evidence recounting. In *CVPR*, pages 2568–2577, 2015.

Model	Top-1 Accuracy (%)	Top-5 Accuracy (%)
RGB	73.0	90.9
Flow	54.5	75.9
Audio	21.6	39.4
VLAD	74.6	90.0
LSTM	73.9	91.6
Multi-stream Sequence Model	76.0	92.7
Fast-forward LSTM	76.3	92.9
Temporal Xception Network	76.1	92.8
Shifting Attention Network	78.7	93.9
Ensemble	81.5	95.6

Table 1. Kinetics validation results.

- [3] S. Hershey, S. Chaudhuri, D. P. W. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold, M. Slaney, R. J. Weiss, and K. Wilson. Cnn architectures for large-scale audio classification. In *arXiv preprint arXiv:1609.09430*, 2017.
- [4] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [5] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [6] F. Li, C. Gan, X. Liu, Y. Bian, X. Long, Y. Li, Z. Li, J. Zhou, and S. Wen. Temporal modeling approaches for large-scale youtube-8m video understanding. *arXiv:1707.04555*, 2017.
- [7] Z. Lin, M. Feng, C. Nogueira dos Santos, M. Yu, B. Xiang, B. Zhou, and Y. Bengio. A Structured Self-attentive Sentence Embedding. *ArXiv e-prints*, Mar. 2017.
- [8] A. Miech, I. Laptev, and J. Sivic. Learnable pooling with Context Gating for video classification. *ArXiv e-prints*, June 2017.
- [9] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014.
- [10] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [11] N. Srivastava, E. Mansimov, and R. Salakhutdinov. Unsupervised learning of video representations using lstms. *ICML*, 2015.
- [12] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *arXiv preprint arXiv:1602.07261*, 2016.
- [13] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. C3D: Generic features for video analysis. *ICCV*, 2015.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention Is All You Need. *ArXiv e-prints*, June 2017.
- [15] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. V. Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016.

Task 2: Trimmed Action Recognition (Kinetics) - TwentyBN's Submission details

Eren Gölge, Raghav Goyal, Valentin Hänel
Twenty Billion Neurons GmbH, Berlin, Germany

eren.golge, raghav.goyal, valentin.haeln @twentybn.com

Abstract

Our approach for action recognition in trimmed single-activity videos exploits both audio and visual streams present in the videos. We train neural networks which capture different statistical structure in the data and finally ensemble the obtained logits to get the final result.

We propose a novel 3D convolutional neural network architecture which employs separable convolution filters for spatial and time domain, additionally we added random dilation in temporal dimension as a regularizer to make the model more robust towards information contained at different time-scales. This single model yielded a validation top-1 accuracy of 70% and is our best single model among others. Our final submission after emsembling the models achieves 78% top-1 accuracy on the validation set

1. Introduction

To recognize activities in trimmed video sequences where each video contains a single activity (label) with duration of not more than 10 secs. We train multiple models and ensemble them to get the final predictions which is shown in the Figure 1

2. Method

2.1. 2D CNN models

Resfeat1 We extract ResNet-101's 2048 dimensional features after the last pool layer for each frame in the video. Further, we pool the features by averaging them and train a MLP on top of it. We experimented with different nonlinearities and obtained the following top-1 accuracies:

- PReLU [5]: 64%
- Maxout units [4] : 65%

Resfeat2 Following the similar architecture as above, but additionally we cluster data into 25 groups by RSOM [3]

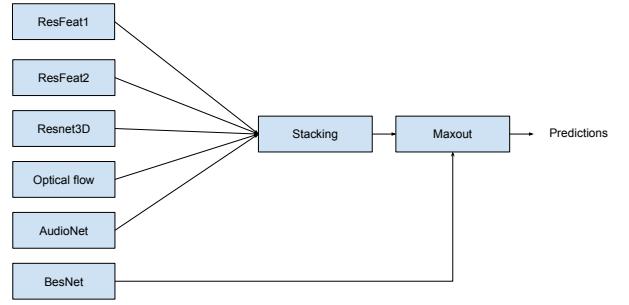


Figure 1. A schematic diagram for our approach

and train a MLP on top of it. We obtained following top-1 accuracies:

- PReLU : 66.2%
- Maxout units : 67.8%

2.2. 3D CNN models

Resnet3D We took the idea from [2] wherein the Inception-v1 filters were inflated in the time domain to obtain ImageNet initialization for 3D CNN model, and inflated Resnet-50 [6]. We obtained a validation accuracy of 64.30% (top-1), 85.58% (top-5)

Optical Flow We took an OpenCV implementation of Gunnar Farneback's algorithm for computing dense optical flow and converted 2-channel optical flow vectors (u, v) into its magnitude and direction and stored them as RGB images for the sake of compression. We used these images in the above described Resnet3D CNN architecture and obtained a validation accuracy of 42.65% (top-1,) 68.09% (top-5)

2.3. AudioNet

We use DeepSpeech-2 architecture [1] to construct audio model and average the obtained RNN hidden state outputs and pass it through an FC layer for the classification task. It solves most of the audio specific classes with high confidence. Overall, we obtained a validation accuracy of 17.86% (top-1), 34.39% (top-5).

2.4. BesNet

This is 3D convolutional network having the same order of layers as Resnet-50. We inflate Resnet-50 trained on ImageNet dataset in order to initialize the network. We decompose temporal and spatial convolutions into 1x1 and 3x3 filters respectively. We use 3D filters only for 1x1 filters and lead them to specialize on temporal relations between channels. Temporal dimension for 3x3 filters are kept as one and let them only to learn spatial relations. One novel property of this architecture is random dilation [7] on temporal layers. Each 1x1 layer randomly picks a dilation factor as 1 or 2. This pertains to learn robust temporal relations between feature maps. Throughout the network, we keep the temporal dimension as it is. At the end, we merge them by max-pooling, following the last spatial average-pooling. It provides a flexible architecture which can be used by any number of frames. We intensify input data augmentation with wide-range scale jitter, random cropping. It allows to regularize the network and also increases the accuracy in case of test-set augmentation. This network did not converge until the deadline but still it gives the best single network performance among all our architectures. It yields 70% (top-1) and 74% (top-1) with test-time augmentation.

Entropy based Hard-mining We apply hard-mining entailing to BesNet training. To this end, we wait the network to get decent results on validation set. Then we start to pick hard instances defined by entropy based metric on the predicted confidences. This leads to pick instances that are not confidently classified as hard instances. We train the network for an additional epoch with these hard instances and we repeat this routine couple of times. It increases top-1 accuracy from 70% to 72% for BesNet. Although hard-mining tend to over-fit, entropy based selections seems to be resilient.

2.5. Ensembling

All our models are ensembled at the end in different methods.

2.5.1 Maxing out

We try different pooling methods. The best performance is obtained by max-pooling. The best combination gives 77%

(top-1) with max-pooling, 71% (top-1) with avg-pooling and 72% (top-1) with majority voting.

2.5.2 Stacking

One alternative technique is stacking. We explain the procedure below. This gives 78% (top-1) on validation set (best submission).

Procedure:

- Train models with normal train and validation split.
- Get class confidence values for validation set from each model.
- Divide validation set into 2 as val1 and val2.
- Train a meta-model on val1 with first-level class confidences.
- Check validation loss on val2

References

- [1] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In *International Conference on Machine Learning*, pages 173–182, 2016.
- [2] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. *arXiv preprint arXiv:1705.07750*, 2017.
- [3] E. Golge and P. Duygulu. Conceptmap: Mining noisy web data for concept learning. In *European Conference on Computer Vision*, pages 439–455. Springer, Cham, 2014.
- [4] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks. *arXiv preprint arXiv:1302.4389*, 2013.
- [5] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [7] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.

Submission to the Kinetics Human Action Recognition Challenge

Saining Xie¹ *

Chen Sun²

Jonathan Huang²

Kevin Murphy²

¹UC San Diego

²Google Research

s9xie@eng.ucsd.edu

{chensun, jonathanhuang, kpmurphy}@google.com

Abstract

This note describes the details of our solution to ActivityNet 2017 Challenge, Kinetics Human Action Recognition track. Motivated by [1], the goal of our submission, is to explore the purely end-to-end trained 3D convolutional neural networks, by inflating modern 2D CNN architectures and fine-tuning from ImageNet initialization. We found that quite surprisingly, using deeper and better (on ImageNet) architectures not necessarily provides better performance for video classification, although their ensemble does. We also find that optical flow based models are still complementary for ensemble, but the gain vanishes.

1. Approach Overview

Our general framework follows the work by Carreira and Zisserman [1], where they propose to "inflate" convolutional neural networks designed for images by replacing $N \times N$ 2D convolutional filters into $N \times N \times N$ 3D filters. The paper also demonstrates the effectiveness of initialization from the 2D counterpart of CNN pretrained on ImageNet by duplicating the weights over time.

Our approach differs from [1] by the following design choices:

- **Base 2D CNN architecture.** We explore several more modern image CNN architectures than Inception-v1 used by [1], they have much higher accuracy on ImageNet but also have more parameters.
- **Initialization strategy.** Rather than replicating ImageNet pre-trained weights over time, we also consider the approach proposed by [2], in which we only set the center channel on the temporal dimension, and zero the others.
- **Ensemble across and within modalities.** The two stream architecture can be seen as a special case of model ensemble where input modalities are different.

We are interested to see if ensemble two RGB networks have the similar effect.

2. Experimental Setup

We start with state-of-the-art CNN models designed for ImageNet classification, including Inception-v1, Inception-v3 and ResNets, whose size and accuracy are shown in Table 1. We inflate the 2D convolution kernels into 3D, and finetune all the parameter layers. For training and evaluation, we average the outputs from the final logits layer over time to generate video-level predictions.

ImageNet model	#Params ($\times 10^6$)	ImageNet Top-1 Acc (%)
Inception-v1	6.8	69.8
Inception-v3	23.2	78.0
ResNet-50	25.5	75.8
ResNet-101	44.3	77.5

Table 1. Number of parameters and ImageNet classification accuracy of our used 2D pretrained models.

Inception-v1. We follow the I3D architecture design proposed in [1]. However, we expand the input frames to 96 instead of 64 in original paper, which results in better performance.

Inception-v3. The 3D conversion is similar to 3D Inception-v3. For the asymmetric filters (e.g. 1×7 and 7×1), we add a singleton temporal dimension and the converted 3D kernels become $1 \times 1 \times 7$ and $1 \times 7 \times 1$, where the first dimension is for time.

ResNet-50 and ResNet-101. We use the standard ResNet-v1 models. Following the configuration of Inception-v1, the first 7×7 convolution kernel is converted to $7 \times 7 \times 7$ with stride 2. The first pooling layer has a kernel shape $1 \times 3 \times 3$ and stride 2 only on spatial dimensions. For the following bottleneck blocks, the kernel of the 1×1 convolution layers are inflated to $1 \times 1 \times 1$ and the 3×3 convolution layers are converted to $1 \times 3 \times 3$.

Initializing 3D filters from 2D Filters. In [1], the 3D filters are constructed by repeating the weights of the 2D filters N times along the temporal dimension, and rescaling them by

*Work done while interning at Google Research.

dividing by N. We find that center initialization [2] (only set the center channel on the temporal dimension and zero others) leads to faster convergence, though the difference on validation error is marginal.

Training Details. We use synchronized SGD with momentum of 0.9 on 64 GPUs. For all experiments, batch size is fixed to 6. Learning rate is 0.1 for Inception-v1 models and 0.01 for Inception-v3 and ResNet models. We train our model for 80,000 steps. These hyper parameters are selected on the validation partition of Kinetics. Optical flow is computed with the standard TV-L1 algorithm and encoded into JPEG images.

3. Key observations

Below we summarize some of our key findings, based on our experiments. We hope to verify these findings more rigorously in the future.

Observation 1: Deeper and better (on ImageNet) models do not necessarily result in better performance after 3D inflation.

Model	Input	Top-1 (%)	Top-5 (%)
Inception-v1	RGB	72.18	90.71
Inception-v1	flow	63.11	84.67
Inception-v3	RGB	70.31	89.18
Inception-v3	flow	61.81	82.87
ResNet-50	RGB	68.32	87.92
ResNet-50	flow	61.23	83.06
ResNet-101	RGB	67.62	87.13

Table 2. Single model classification accuracy on Kinetics validation set. The accuracy is evaluated on a single center crop of 224×224 pixels.

It is generally believed that better ImageNet classification results (i.e. better “backbone architectures”), typically can transfer to other tasks in different domains, e.g. object detection [3]. However, as can be observed from Table 2, the current approach is not able to fully utilize the model capacity and representation power of Inception-v3 and ResNet models. This indicates that more careful architecture design may be needed than simply inflating 2D models over time.

Observation 2: Ensembling can help somewhat, especially when we combine RGB and flow.

In Table 3 we show the benefits of combining multiple Inception-v1 models, to illustrate the relative gain of different choice of ensembles. We see that ensembling two RGB

Model	Top-1 (%)	Top-5 (%)
1 RGB _{Inception-v1}	73.21	91.42
2 RGB _{Inception-v1}	74.43($\uparrow 1.22\%$)	91.88($\uparrow 0.46\%$)
1 RGB _{Inception-v1} + 1 flow _{Inception-v1}	75.05($\uparrow 1.84\%$)	92.08($\uparrow 0.66\%$)
1 RGB _{Inception-v1} + 1 RGB _{Inception-v3}	74.96($\uparrow 1.75\%$)	92.23($\uparrow 0.81\%$)

Table 3. Inception-v1 (multi-scale) model ensemble results on Kinetics Validation Set.

inception-v1 models helps a bit, but that ensembling one RGB and flow helps a bit more. Interestingly, ensembling one RGB inception-v1 with one RGB inception-v3 almost matches the benefits of combining RGB and flow. Our final solution is an ensemble of RGB and flow models trained with Inception-v1, Inception-v3, ResNet-50 and ResNet-101. With this, we achieve a final accuracy of 76.97% (top-1) and 93.04% (top-5) on validation set.

Observation 3: Using wider temporal kernels (more input frames) is helpful.

For our Inception-v1 model, using $1.5 \times$ more input frames (from 64 frames to 96 frames) gives better performance.

Observation 4: Using multi-scale testing is helpful.

For 2D images, evaluating at multiple scales often helps. In Table 4, we show a similar benefit on Kinetic video data, especially for RGB models. (For crop sizes larger than 224, we apply the model fully-convolutionally, and then spatially average the predictions after the softmax layer.) Flow models are more scale sensitive, and the improvement with multi-scale testing is minor.

Observation 5: Inflating more convolution layers to 3D seems always helpful.

For different architectures, we find that inflating more 2D filters into 3D (therefore larger temporal receptive field) leads to better performance on both training and validation set. For example, a naive 2D ResNet-50-RGB model applied on 3D videos, where all 3×3 filters are converted into $1 \times 3 \times 3$ gives a top-1 accuracy of 61.0%, whereas after using $3 \times 3 \times 3$ inflated filters on all bottleneck blocks, boosts this to 68.32%. Generally speaking, we find that inflating more convolution filters in more layers leads to increasingly better accuracy.

Model	Scale	Input	Top-1 (%)	Top-5 (%)
Inception-v1	224	RGB	72.18	90.71
Inception-v1	256	RGB	72.72	91.01
Inception-v1	320	RGB	71.72	90.52
Inception-v1	224 + 256 + 320	RGB	73.21	91.42
Inception-v3	224 + 256 + 320	RGB	70.97	89.64
ResNet 50	224 + 256 + 320	RGB	69.57	88.55
ResNet 101	224 + 256 + 320	RGB	68.55	87.91
Inception-v1	224	flow	63.11	84.67
Inception-v1	256	flow	63.22	84.70
Inception-v1	320	flow	59.79	82.14
Inception-v1	224 + 256 + 320	flow	63.38	84.79
Inception-v1	224 + 256	flow	63.56	84.70

Table 4. Multi-scale testing results.

4. Acknowledgement

We would like to thank the authors of the Kinetics dataset for their kind support; and Joao Carreira and Brian Zhang for helpful discussions.

References

- [1] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. *CVPR*, 2017. [1](#)
- [2] C. Feichtenhofer, A. Pinz, and R. Wildes. Spatiotemporal multiplier networks for video action recognition. In *CVPR*, 2017. [1](#), [2](#)
- [3] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, et al. Speed/accuracy trade-offs for modern convolutional object detectors. *CVPR*, 2017. [2](#)

ActivityNet 2017: Improved Benchmarks Using 3-D Convolutional Networks for Trimmed Video Sequences.

Brendan Duke¹, Fabien Baradel², Thorsteinn H. Jonsson¹, Terrance Devries¹, Christian Wolf², Graham Taylor¹, and Julien Mille³

¹School of Engineering, University of Guelph, Guelph, ON, Canada

²Univ Lyon, INSA-Lyon, CNRS, LIRIS, F-69621, Villeurbanne, France

³Laboratoire d’Informatique de l’Université de Tours (EA 6300), INSA Centre Val de Loire, 41034 Blois, France

{bduke, tjonsson, tdevries, gtaylor}@uoguelph.ca, {fabien.baradel, christian.wolf}@liris.cnrs.fr, julien.mille@insa-cvl.fr

Abstract

This technical report details the methods used by the University of Guelph - LIRIS team for the ActivityNet Trimmed Action Recognition - 2017 challenge. We present improved benchmarks with two streams of three dimensional convolutional networks (3D-CNN) and show results for different strategies for aggregating spatio-temporal features. We compare the difference between considering a Long Short-Term Memory Network and a 1D Convolutional Network for learning high level temporal dependencies. We also establish benchmarks for RGB and “depth and motion” datastreams using transfer learning techniques. Using a single model, our best method improves current benchmarks by Kay et al.[7] by 3.65 points using the averaged top-1 and top-5 error metric. Our final submission is an ensemble of six models: five based on RGB and one based on depth and motion, achieving an averaged top-1 and top-5 error percentage of 24.3% on the test set.

1. Introduction

Learning spatio-temporal structure from videos has remained a challenging task in the computer vision community, mainly due to the computational resources needed for the task, but also due to the lack of large and robust datasets. Work on learned temporal pooling with LSTMs of 3D-CNN features can already be seen back in 2011 with the work of Baccouche et al. [1] While recent approaches [6, 12, 14] reached interesting results on standard benchmark datasets HMDB-51 [9], UCF-101 [13] they show their weakness on new large scale datasets [4, 7]. Meanwhile, convolutional neural networks (CNN) have shown human-level perfor-

mance [11] in image classification and object recognition. Carreira et al. [2] show promising results by inflating 2D CNN model weights pre-trained on ImageNet to 3D CNN weights on the new Kinetics dataset. Their method makes the learning of spatio-temporal features benefit from a powerful model already trained at extracting useful spatial features. Our model extends this idea to variants of the inflated two-stream model as well as different strategies for pooling features across the temporal dimension.

2. Methodology

Our investigation focuses on improvements that can be made over the established benchmarks with 3D convolutional networks. Transfer learning from 2D convolutional networks to 3D convolutional networks is done via the inflation scheme introduced by [7]. For a given video sequence we perform global average pooling on short segments of consecutive frames. Given the log probabilities obtained from the independent streams we learn higher-level temporal features by aggregating using either a 1D convolutional neural network or a LSTM neural network. Similar model comparisons can be seen in literature on character based language models [8, 5, 16].

Our model consists of two streams of 3D-inflated residual networks, which are obtained by performing transfer learning from ImageNet via CNN as well as from the DeMoN-Network structure in motion data [15]. Each stream is trained separately and ensembled for the final prediction.

In all experiments below, frames were sampled from videos at 25 frames per second. Videos were resampled to 256×256 , and then randomly cropped down to 224×224 for input to the ResNet-3D models.

Standard data augmentation strategies were used.

2.1. Inflated 3D CNN

We follow the work done in [2] to inflate 3D CNN kernels from 2D CNN weights, which is to expand the kernels of pretrained architectures by tiling them across the temporal dimension.

This allows us to extract spatio-temporal features from video segments, the temporal aspect of which would be uncaptured by 2D CNNs.

We select the powerful ResNet architectures (ResNet-50 and ResNet-101) and inflated them into 3D ResNet. Inflating an architecture affects the filters and pooling kernels. We endowed them by adding a temporal dimension (e.g. a convolutional filter of size $N \times N$ becomes a filter of size $N \times N \times N$), following the study of filter dimensions in [14], which found $3 \times 3 \times 3$ kernels advantageous over various other filter shapes.

The inflated 3D ResNet gives a spatio-temporal feature vector of size 2048 from which the classification predictions are performed.

We train the RGB ResNet-3D model using 32-frame snippets. During the evaluation we average the predictions from 14 different snippets.

2.2. LSTM and 1D-Conv

The main advantage of using 3D-Convolutional network is that spatio-temporal features are learned in harmony. However, since gradients are not being propagated across different time segments in sequence, the model lacks the ability to capture long range dependencies. To test this hypothesis we employ an LSTM and a 1D-CNN at the top of the 3D-CNN.

In all cases the inputs to our model are sequences of 128 RGB frames from a video that is decomposed into 7 sub-sequences of 32 frames, each overlapping every 16 frames. We use the inflated 3D-CNN described above to extract a feature vector from each sub-sequence via 3D global average pooling. For some given decoder, the sequence of feature vectors is passed through it and classification is subsequently done with a soft-max. We set the LSTM with a hidden state of 1024 neurons.

As an additional method of learned temporal pooling of features, besides LSTM, 1D convolutions are used analogous to the use of 1D convolutions for sentence classification in [16].

Sets of filters of different region sizes 7, 5, 3 and 1 are used and the input sequence is not padded, such that these filter sets produce outputs of lengths 1, 3, 5 and 7 from the input feature vector sequences, respectively. The entire 2048×7 sequence of feature vectors extracted by the ResNet-3D from 32-frame snippets is convolved by each of these filters, which each produce outputs of 512 channels.

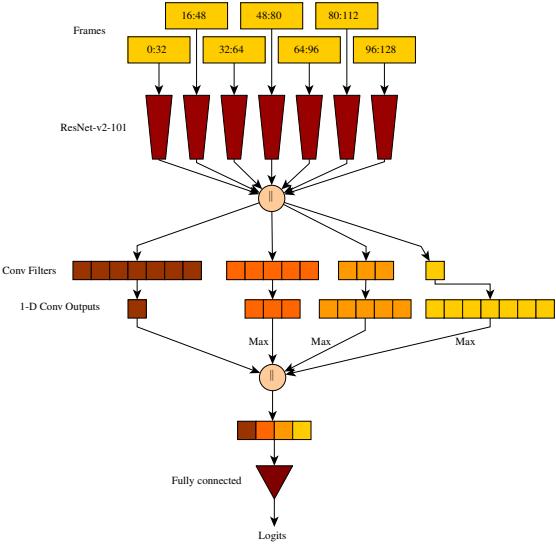


Figure 1: Linear convolution at the top of the inflated 3D CNN. Seven 32-frame snippets, overlapping by 16-frames each for a total receptive field of 128 frames, are input to ResNet-v2-101 feature extractors. The \parallel operator denotes concatenation. The extracted features are concatenated into a sequence, which is convolved with 1D filters of sizes 1, 3, 5 and 7, the result of which is 1-max pooled, concatenated into a 2048-channel vector, and input to a fully-connected prediction layer.

The 512-output channel outputs from the 1D convolutions go through 1-max pooling and are concatenated channel-wise, to produce a single vector of 2048 channels, which is input to a fully-connected layer.

The 1D convolutional model architecture is illustrated in Figure 1.

Weights of the 3D CNN are initialized with weights from applying classification over the 3D CNN only (previous section). For evaluation, we extract 3 different 128-frame snippets and average predictions.

2.3. Depth & Motion

We use a DeMoN model, pre-trained from the code released along with [15], to incorporate depth, surface normal and optical flow information into the “depth and motion” stream of the activity classifier.

As illustrated in Figure 2, depth, surface normal, optical flow and optical flow confidence outputs from the DeMoN bootstrap network are concatenated as an input to the depth and motion stream ResNet-v2-50 with dimensions as follows: 8-channels, width and height of 56 pixels, and temporal length of 16-frames.

To produce optical flow and depth predictions, image pairs 20 frames apart are passed to the DeMoN model.

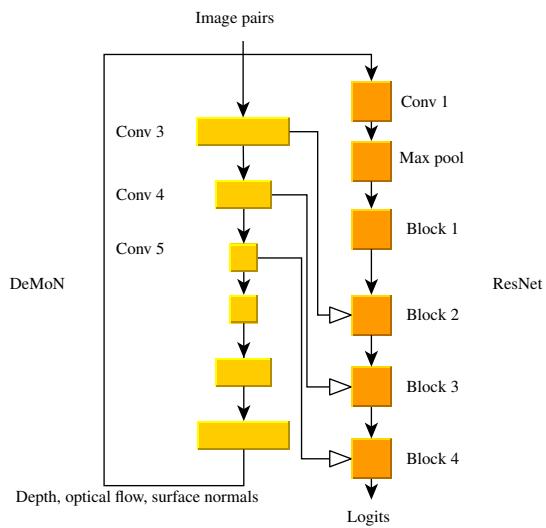


Figure 2: Depth and motion stream ResNet. Output predictions from the pre-trained DeMoN model of [15] are input to the depth and motion stream. Furthermore, feature maps from the DeMoN model are “injected” into the depth and motion ResNet in ResNet blocks 2, 3 and 4.

In addition to the raw output predictions from the DeMoN bootstrap network, we also “inject” features from layers `conv3_1`, `conv4_1` and `conv5_1` into the depth and motion stream ResNet at the layers where the number of channels of the respective optical flow and depth estimation DeMoN models are equal to the number of channels of the ResNet.

In order to reduce the temporal length of the features injected from the DeMoN model, those features are input to a sequence of n ResNet bottleneck-v2 units, where the temporal length of the ResNet feature maps at the point that the DeMoN features will be injected is $16/2^n$. Each of these “injection units” has a stride of two in the temporal dimension, and preserves spatial dimensions. The resulting feature maps are added, before the filter-sized $3 \times 3 \times 3$ convolution, to the reduced-dimensionality feature maps in bottleneck units of the depth and motion stream ResNet.

3. Experiments

We report results for the Kinetics dataset on the validation set in Table 1.

We trained a ResNet-v2-50 3D from Xavier initialization, i.e. without pre-training, on the Kinetics dataset and achieved a top-5 accuracy of 83.0% after 110 epochs. Achieving similar accuracy to the inflated ResNet-50 3D shows that it is possible to use Kinetics to train 3D convolutional models from scratch, without first pre-training on

ImageNet.

The advantage of deeper networks in activity recognition in video is shown by the improved classification accuracy of the 3D ResNet-101 model, which improved upon the averaged top-1 and top-5 error of the 3D ResNet-50 model by 1.6 percentage points.

The LSTM and 1D convolutional models showed an improvement of 0.3 and 0.4 percentage points in top-5 error compared with averaging ResNet-101 predictions over 14 uniformly sampled 32-frame snippets. However, these “temporal pooling” models performed worse in top-1 error compared to uniform ResNet-101 predictions by 0.7 (LSTM) and 0.9 (1D convolutional) percentage points.

The depth and motion stream trained on DeMoN predictions and feature maps achieved an averaged top-1 and top-5 error of 38.3%, improving over the flow-only baseline of [7] by 1.0 percentage points.

4. Implementation

Our codebase is written in TensorFlow. To work with videos in a sensible manner we implemented a tensorflow operation that can decompress H264 files. This allows us to train our network on a compressed dataset of about 75 GB.

Based on the code provided by the TensorFlow authors at [10], we implemented a `StagingArea` pipeline that decouples variable weight updates from gradient computation. The implementation is able to train a ResNet-50 3D for 200K iterations in 40 minutes, on eight 12GB NVIDIA P100 GPUs distributed over two compute nodes.

A per-GPU batch size of 10 examples is used, for a total batch size of 80 examples. We used the RMSProp optimizer in all experiments. For the experiments that started with inflated pre-trained ImageNet weights, we used an initial learning rate of 0.004 and an exponential learning rate decay that decayed by a factor of 0.16 every 20 epochs for 90 epochs total. For training the ResNet-v2-50 3D from scratch, a staircase learning rate starting at 0.1 was used, multiplying by a factor of 0.16 every 30 epochs for the first 100 epochs, after which an exponential decay was used for 30 epochs.

Code will be released at a later date.

5. Conclusion

We presented a methodology based on inflated ResNet-3D models that achieved an average top-1 and top-5 error of 23.4% on the test set of the ActivityNet Trimmed Action Recognition challenge using an ensemble of six models. We did not find significant improvement of accuracy using either 1D convolution or LSTM learned pooling strategies compared with averaging predictions on 32-frame snippets over the entire 10s video clips. We improved over baseline

Methods	Top-1	Top-5	Average Error
3D ResNet-50 Xavier Initialization [3]	60.7	83.0	28.2
Inflated 3D ResNet-50	62.5	83.9	26.8
Inflated 3D ResNet-101	64.6	85.1	25.2
Inflated 3D ResNet-101 + LSTM (128 frames)	63.9	85.4	25.3
Inflated 3D ResNet-101 + 1D-Conv (128 frames)	63.7	85.5	25.4
Depth&Motion	49.5	74.0	38.3

Table 1: Results on Kinetics dataset - validation set (accuracies in %)

benchmarks in [7] by an averaged top-1 and top-5 error percentage of 3.65 points.

References

- [1] M. Baccouche, F. Mamalet, C. Wolf, C. Garcia, and A. Baskurt. Sequential deep learning for human action recognition. In *Proceedings of the Second International Conference on Human Behavior Understanding*, HBU'11, pages 29–39, Berlin, Heidelberg, 2011. Springer-Verlag. [1](#)
- [2] J. Carreira and A. Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. *CVPR*, 2017. [1, 2](#)
- [3] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS'10)*. Society for Artificial Intelligence and Statistics, 2010. [4](#)
- [4] R. Goyal, S. E. Kahou, V. Michalski, J. Materzyska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, F. Hoppe, C. Thurau, I. Bax, and R. Memisevic. The “something something” video database for learning and evaluating visual common sense. [1](#)
- [5] R. Józefowicz, O. Vinyals, M. Schuster, N. Shazeer, and Y. Wu. Exploring the limits of language modeling. *CoRR*, abs/1602.02410, 2016. [1](#)
- [6] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014. [1](#)
- [7] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman. The Kinetics Human Action Video Dataset. 2017. [1, 3, 4](#)
- [8] Y. Kim, Y. Jernite, D. Sontag, and A. M. Rush. Character-aware neural language models. *CoRR*, abs/1508.06615, 2015. [1](#)
- [9] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: a large video database for human motion recognition. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011. [1](#)
- [10] The TensorFlow Authors. Performance benchmarks, 2017. Code available at <https://github.com/tensorflow/benchmarks>. [3](#)
- [11] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014. [1](#)
- [12] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. *NIPS*, 2014. [1](#)
- [13] K. Soomro, A. R. Zamir, and M. Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *CoRR*, abs/1212.0402, 2012. [1](#)
- [14] D. Tran, L. D. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 4489–4497, 2015. [1, 2](#)
- [15] B. Ummenhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, and T. Brox. Demon: Depth and motion network for learning monocular stereo. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [1, 2, 3](#)
- [16] Y. Zhang and B. C. Wallace. A sensitivity analysis of (and practitioners’ guide to) convolutional neural networks for sentence classification. *CoRR*, abs/1510.03820, 2015. [1, 2](#)

MSR Asia MSM at ActivityNet Challenge 2017: Trimmed Action Recognition, Temporal Action Proposals and Dense-Captioning Events in Videos

Ting Yao, Yehao Li, Zhaofan Qiu, Fuchen Long, Yingwei Pan, Dong Li, and Tao Mei

Microsoft Research, Beijing, China

{tiyao, tmei}@microsoft.com

Abstract

This notebook paper presents an overview and comparative analysis of our systems designed for the following three tasks in ActivityNet Challenge 2017: trimmed action recognition, temporal action proposals and dense-captioning events in videos.

Trimmed Action Recognition (TAR): We investigate and exploit multiple spatio-temporal clues for trimmed action recognition (TAR) task, i.e., frame, short video clip and motion (optical flow) by leveraging 2D or 3D convolutional neural networks (CNNs). The mechanism of different quantization methods is studied as well. Furthermore, improved dense trajectory with fisher vector encoding over the whole trimmed video is utilized. All activities are finally classified by late fusing the predictions of one-versus-rest linear SVMs learnt on each clue.

Temporal Action Proposals (TAP): To generate temporal action proposals from videos, a three-stage workflow is particularly devised for TAP task. Given an untrimmed video, our system firstly generates an actionness curve via a snippet-level actionness classifier. The temporal actionness grouping scheme is then exploited over actionness curve to produce proposal candidates. Finally, a proposal re-ranking procedure is incorporated to select high-quality proposals via a proposal-level actionness classifier.

Dense-Captioning Events in Videos (DCEV): For DCEV task, we firstly adopt our temporal action proposal system mentioned above to localize temporal proposals of interest in video, and then generate the descriptions for each proposal. Specifically, RNNs encode a given video and its detected attributes into a fixed dimensional vector, and then decode it to the target output sentence. Moreover, we extend the attributes-based CNNs plus RNNs model with policy gradient optimization and retrieval mechanism to further boost video captioning performance.

1. Introduction

Recognizing activities in videos is a challenging task as video is an information-intensive media with complex variations. In particular, an activity may be represented by different clues including frame, short video clip, motion (optical flow) and long video clip. In this work, we aim at investigating these multiple clues to activity classification in trimmed videos, which consist of a diverse range of human focused actions. However, most of the natural videos in the real world are untrimmed videos with complex activities and unrelated background/context information, making it hard to directly recognize activities in them. One possible solution is to quickly localize temporal chunks in untrimmed videos containing human activities of interest and then conduct activity recognition over these temporal chunks, which largely simplifies the activity recognition for untrimmed videos. Generating such temporal action chunks in untrimmed videos is known as the task of temporal action proposals, which is also exploited in this work.

In addition to the above two tasks tailored to activity which is usually the name of action/event in videos, the task of dense-captioning events in videos is explored here which goes beyond activities by describing numerous events within untrimmed videos with multiple natural sentences.

The remaining sections are organized as follows. Section 2 presents all the features which will be adopted in our systems, while Section 3 details the feature quantization strategies. Then the descriptions and empirical evaluations of our systems for three tasks are provided in Section 4-6 respectively, followed by the conclusions in Section 7.

2. Video Representations

We extract the video representations from multiple clues including frame, short clip, motion and long clip.

Frame. To extract frame-level representations from video, we uniformly sample 25 frames for each video/proposal, and then use pre-trained 2D CNNs as

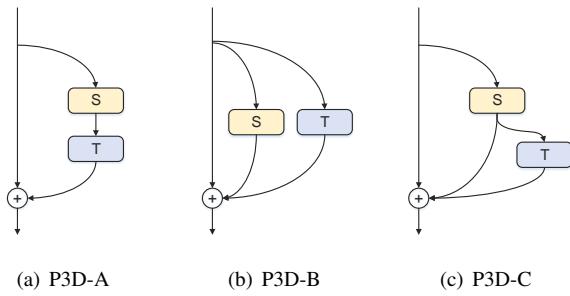


Figure 1. Three Pseudo-3D blocks.

frame-level feature extractors. We choose the most popular 2D CNNs in image classification—ResNet [4].

Short Clip. In addition to frame, we take the inspiration from the most popular 3D CNN architecture C3D [18] and devise a novel Pseudo-3D Residual Net (P3D ResNet) architecture [15] to learn spatio-temporal video clip representation in deep networks. Particularly, we develop variants of bottleneck building blocks to combine 2D spatial and 1D temporal convolutions, as shown in Figure 1. The whole P3D ResNet is then constructed by integrating Pseudo-3D blocks into a residual learning framework at different placements. Our P3D ResNet model is pre-trained on Sports-1M dataset [5]. We fix 16 frames as the length of short clip, and sample rate is set to 25 per video.

Motion. To model the change of consecutive frames, we apply another CNNs to optical flow “image,” which can extract motion features between consecutive frames. When extracting motion features, we follow the setting of [21], which fed 32 optical flow images, consisting of two-direction optical flow from 16 consecutive frames, into ResNet/P3D ResNet network in each iteration. The sample rate is also set to 25 per video.

Long Clip. For long/trimmed clip, we choose the state-of-the-art hand-crafted features—improved dense trajectory (iDT) [20] on each trimmed clip. Specifically, trajectory feature, histogram of oriented gradients (HOG), histogram of flow (HOF), and motion boundary histogram (MBH) are computed for each trajectory obtained by tracking points in video clips. Furthermore, Fisher vector encoding is used to quantize the features and create high dimensional representations for each clip.

3. Feature Quantization

In this section, we describe two quantization methods to generate video-level representations from frame-level or clip-level features.

Average Pooling. Average pooling is the most common method to extract video-level features from consecutive frames, short clips and long clips. For a set of frame-level or clip-level features $F = \{f_1, f_2, \dots, f_N\}$, the video-

level representations are produced by simply averaging all the features in the set:

$$R_{pooling} = \frac{1}{N} \sum_{i:f_i \in F} f_i , \quad (1)$$

where $R_{pooling}$ denotes the final representations.

Deep Quantization. Moreover, we present our recently proposed network-based quantization method called Deep Quantization (DQ) [14]. A generative neural network with parameters θ is trained on the top of feature extraction network. Then, following the fisher kernel method, the video-level representations are defined as

$$\begin{aligned} L_{Generative}(\theta) &= \sum_{f \in TrainingSet} -\log p(f, \theta) \\ \hat{\theta} &= \arg \max_{\theta} L_{Generative}(\theta) \\ R_{DQ} &= \text{normalize}\left(\sum_{i:f_i \in F} \frac{\partial(-\log p(f_i, \hat{\theta}))}{\partial \hat{\theta}}\right) \end{aligned} , \quad (2)$$

where $p(f, \theta)$ is the generative network output. After optimizing parameters θ , the gradient calculating and accumulating can be processed end-to-end during backpropagation, no extra storage is required. To further improve the ability of representations, we propose a semi-supervised optimizing function as:

$$\begin{aligned} L(\theta) &= L_{Generative}(\theta) + \lambda L_{Classification}(\theta) \\ \hat{\theta} &= \arg \max_{\theta} L(\theta) \\ R_{DQ} &= \text{normalize}\left(\sum_{i:f_i \in F} \frac{\partial(-\log p(f_i, \hat{\theta}))}{\partial \hat{\theta}}\right) \end{aligned} . \quad (3)$$

Readers can refer to [14] for more technical details of our deep quantization network.

4. Trimmed Action Recognition

4.1. System

Our trimmed action recognition framework is shown in Figure 2 (a). In general, the trimmed action recognition process is composed of three stages, i.e., multi-stream feature extraction, feature quantization and prediction generation. For deep feature extraction, we follow the multi-stream approaches in [8, 13], which represented input video by a hierarchical structure including individual frame, short clip and consecutive frame. In addition to deep features, one most complementary hand-crafted feature, i.e., iDT, is exploited to further enrich the video representations. After extraction of raw features, different quantization and pooling methods are utilized on different features to produce global representations of each trimmed video. Finally, a linear SVM is trained on each kind of video representations and the predictions from multiple SVMs are combined by linearly fusion. When training SVM, we fix $C = 100$.

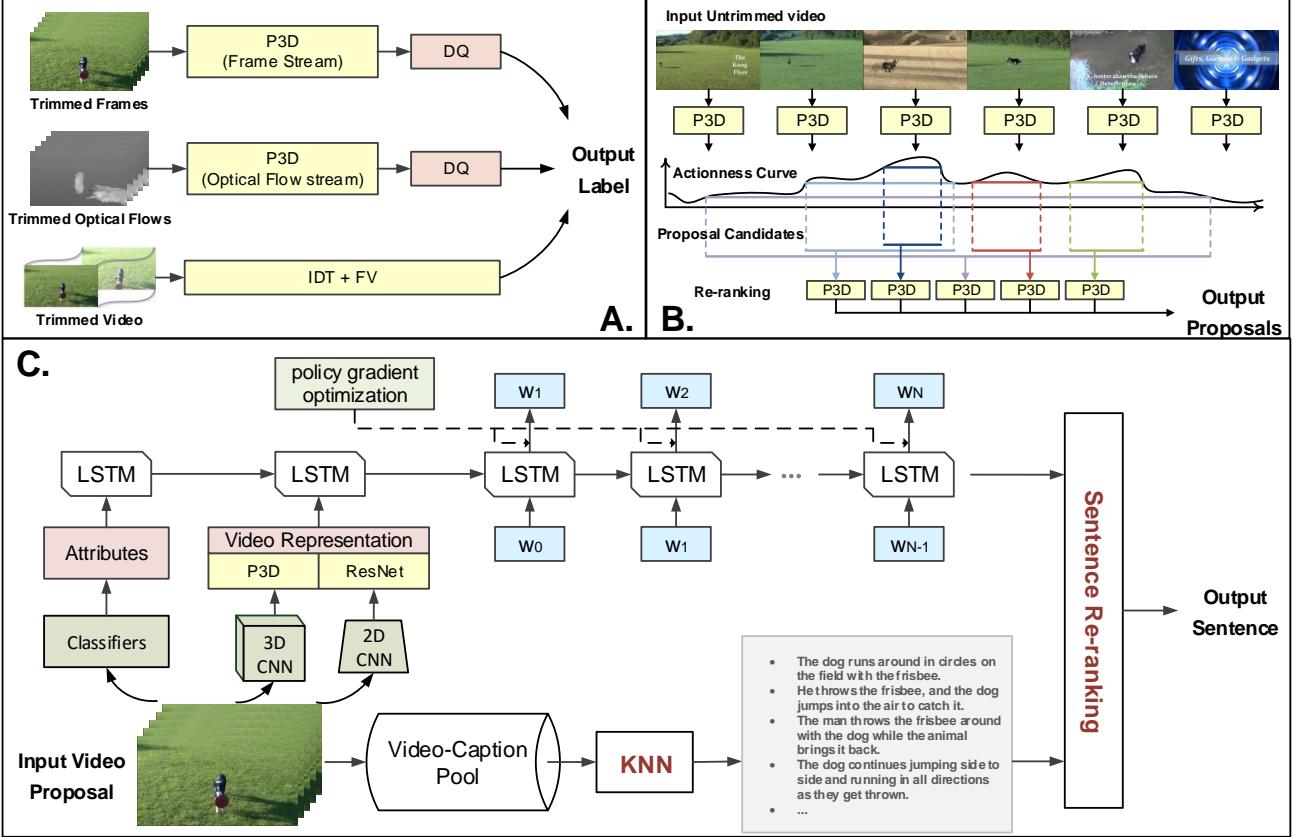


Figure 2. Frameworks of our proposed (a) trimmed action recognition system, (b) temporal action proposals system and (c) dense-captioning system.

4.2. Experiment Results

Table 1 shows the performances of all the components in our trimmed action recognition system. Overall, our Deep Quantization on P3D ResNet achieves the highest top1 accuracy (72.66%) and top5 accuracy (90.74%) of single component. For the final submission, we train the SVMs using training and validation sets. All the components are linearly fused using the weights tuned on validation set.

5. Temporal Action Proposals

5.1. System

Figure 2 (b) shows the framework of temporal action proposals, which is mainly composed of three stages:

Actionness curve generation. We treat every 16 continuous frames as one snippet and the stride size is 8 frames. Then, similar to video highlight detector in [22], a binary actionness classifier is trained over snippets to distinguish whether the snippets contain human activities. Accordingly, an actionness curve can be generated by accumulating all the actionness probabilities of snippets via snippet-level actionness classifier.

Temporal actionness grouping. Given an actionness curve, the classic watershed algorithm [17] is utilized to produce a set of “basins” corresponding to the temporal region with high actionness probability. Then, the temporal actionness grouping scheme [25] is leveraged to connect small basins, resulting in proposal candidates. Finally, the highly overlapped proposal candidates are filtered out via Non-maximal suppression.

Proposal re-ranking. To select the action proposals with high actionness probabilities, we additionally train the proposal-level actionness classifier to measure the actionness probability of each proposal candidate and then re-rank all the proposal candidates. In our experiments, only the top 100 proposals are finally outputted.

5.2. Experiment Results

Table 2 shows the results of actionness classifiers trained with different 2D/3D architectures (i.e., ResNet [4] and P3D ResNet [15]). Each 2D/3D architecture is pre-trained on different sources (e.g., ImageNet [2], Sports1M [5] and Kinetics [6]). For all the single stream runs w or w/o re-ranking scheme, the setting based on P3D ResNet pre-trained on Kinetics achieves the highest AUC. Moreover, by

Table 1. Comparison of different components in our framework on Kinetics validation set for trimmed action recognition task.

Stream	Feature	Layer	Quantization	Top1	Top5
Frame	ResNet	pool5	Ave	70.70%	89.75%
	ResNet	res5c	DQ	71.50%	90.20%
Short Clip	P3D ResNet	pool5	Ave	71.24%	90.01%
	P3D ResNet	res5c	DQ	72.66%	90.74%
Long Clip	iDT+FV	-	-	45.09%	69.73%
Motion	ResNet	pool5	Ave	59.84%	82.54%
	ResNet	res5c	DQ	61.03%	83.51%
	P3D ResNet	pool5	Ave	61.92%	84.19%
	P3D ResNet	res5c	DQ	63.24%	85.53%

Table 2. Area Under the average recall vs. average number of proposals per video Curve (AUC) of different 2D/3D architectures and pre-trained sources on ActivityNet validation set for temporal action proposals task.

Network	Pre-trained	Re-ranking	AUC
ResNet	ImageNet		56.96%
ResNet	Kinetics		59.75%
P3D ResNet	Sports1M		58.79%
P3D ResNet	Kinetics		59.90%
ResNet	ImageNet	✓	59.03%
ResNet	Kinetics	✓	60.13%
P3D ResNet	Sports1M	✓	60.76%
P3D ResNet	Kinetics	✓	61.13%
Fusion all			63.12%

additionally incorporating the re-ranking scheme, our system is consistently improved under different deep architectures. For the final submission, we fusion all the proposals from the eight streams with different settings and then select the top 100 proposals based on their weighted actionness probabilities. The linear fusion weights are tuned on validation set.

6. Dense-Captioning Events in Videos

6.1. System

The main goal of dense-captioning events in videos is jointly localizes temporal proposals of interest in videos and then generate the descriptions for each proposal/video clip. Hence we firstly leverage the temporal action proposal system described above in Section 5 to localize temporal proposals of events in videos (50 proposals for each video). Then, given each temporal proposal (i.e., video segment describing one event), our dense-captioning system runs two different video captioning modules in parallel—the generative module for generating caption via the LSTM-based sequence learning model, and the retrieval module which can directly copy sentences from other visually similar video segments through KNN. Finally, a sentence re-ranking module is exploited to rank and select the final most

consensus caption from the two parallel video captioning modules by considering the lexical similarity among all the sentence candidates. The overall architecture of our dense-captioning system is shown in Figure 2 (c).

Generative module with LSTM. Taking inspiration from the recent successes of probabilistic sequence models leveraged in image/video captioning [9, 10, 11, 19, 23], we follow our previous state-of-the-art image captioning model [24] and formulate the generative video captioning module in an end-to-end fashion based on LSTM which encodes the given video segment and its detected attributes/categories into a fixed dimensional vector and then decodes it to the target output sentence. Specifically, the third design LSTM-A₃ in [24] which firstly encodes attribute representations into LSTM and then transforms video representations into LSTM at the second time step is adopted as the basic architecture. Here, we uniform sample 2 frames/clips per second for each video segment and each word in the sentence is represented as “one-hot” vector (binary index vector in a vocabulary). For the input video representations, we take the output of 2048-way *pool5* layer from the ResNet [4] pre-trained on Kinetics dataset [6] and 2048-way *pool5* layer from P3D ResNet [15] pre-trained on Sports-1M video dataset [5] as frame/clip representation respectively, and then concatenate the features from ResNet and P3D ResNet as the input video representation. For representation of attributes/categories, we treat all the 200 categories on Activitynet dataset [1] as the high-level semantic attributes and train the attribute detectors with our previous video classification system [12], resulting in the final 200-way vector of probabilities. The dimension of the input and hidden layers in LSTM are both set to 1,024.

Furthermore, different from the common training strategy with maximum likelihood estimation (MLE) in LSTM-A₃, we employ the policy gradient optimization method with reinforcement learning [16] to boost the video captioning performances specific to both CIDEr-D and METEOR metrics. Moreover, it should be noted that we additionally incorporate context information from other neighboring events into this generative module like [7].

Table 3. Performance of our proposed dense-captioning models on ActivityNet captions validation set, where B@N, M, R and C are short for BLEU@N, METEOR, ROUGE-L and CIDEr-D scores. All values are reported as percentage (%).

Model	B@1	B@2	B@3	B@4	M	R	C
LSTM-A ₃	17.50	9.62	5.54	3.38	7.71	13.27	16.08
LSTM-A ₃ + policy gradient	17.49	9.73	5.38	3.07	8.47	14.28	13.82
LSTM-A ₃ + policy gradient + retrieval	17.27	9.70	5.39	3.13	8.73	14.29	14.75

Retrieval module with KNN. Another direction of image/video captioning is search-based approaches which “generate” sentence for an image/video by directly copying sentences from other visually similar images/videos. Although the approaches in this dimension cannot produce novel descriptions, it indeed can achieve human-level descriptions as all sentences are from existing human-generated sentences. Hence we design the retrieval module in this dimension to leverage the “crowdsourcing” human intelligence for producing diverse sentences from other angles. In particular, we utilize KNN to find the visually similar video segments based on the extracted video representations. The captions associated with the top similar video segments are regarded as sentence candidates in retrieval module. In the experiment, we mainly choose the top 300 nearest neighbors for generating sentence candidates.

Sentence re-ranking. Given the sentence candidates generated by generative and retrieval modules for input video segment, we need to re-rank all the sentence candidates and select the best one as the final output result. Inspired by [3], we treat the consensus sentence which has the highest average lexical similarity to the other candidates as the best one. Specifically, we linearly fuse two kinds of sentence similarities (i.e., CIDEr-D and METEOR) as the lexical similarity between two sentence candidates.

6.2. Experiment Results

Table 3 shows the performances of our proposed dense-captioning models. Here we compare three variants derived from our proposed dense-captioning framework. In particular, by additionally incorporating the policy gradient optimization scheme into the basic LSTM-A₃ architecture, we can clearly observe the performance boost in METEOR. Moreover, our dense-captioning model (LSTM-A₃ + policy gradient + retrieval) is further improved by injecting the sentence candidates from retrieval module in METEOR.

7. Conclusion

In ActivityNet Challenge 2017, we mainly focused on multiple visual features, different strategies of feature quantization and video captioning from different dimensions. Our future works include more in-depth studies of how fusion weights of different clues could be determined to boost the action recognition/temporal action proposals performance and how to generate open-vocabulary sentences for events in videos.

References

- [1] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, 2015.
- [2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [3] J. Devlin, S. Gupta, R. Girshick, M. Mitchell, and C. L. Zitnick. Exploring nearest neighbor approaches for image captioning. *arXiv preprint arXiv:1505.04467*, 2015.
- [4] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [5] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [6] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [7] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, and J. C. Niebles. Dense-captioning events in videos. *arXiv preprint arXiv:1705.00754*, 2017.
- [8] Q. Li, Z. Qiu, T. Yao, T. Mei, Y. Rui, and J. Lu-o. Action recognition by learning deep multi-granular spatio-temporal video representation. In *ICMR*, 2016.
- [9] Y. Pan, T. Mei, T. Yao, H. Li, and Y. Rui. Jointly modeling embedding and translation to bridge video and language. In *CVPR*, 2016.
- [10] Y. Pan, Z. Qiu, T. Yao, H. Li, and T. Mei. Seeing bot. In *SIGIR*, 2017.
- [11] Y. Pan, T. Yao, H. Li, and T. Mei. Video captioning with transferred semantic attributes. In *CVPR*, 2017.
- [12] Z. Qiu, D. Li, C. Gan, T. Yao, T. Mei, and Y. Rui. Msr asia msm at activitynet challenge 2016. In *CVPR workshop*, 2016.
- [13] Z. Qiu, Q. Li, T. Yao, T. Mei, and Y. Rui. Msr asia msm at thumos challenge 2015. In *THUMOS’15 Action Recognition Challenge*, 2015.

- [14] Z. Qiu, T. Yao, and T. Mei. Deep quantization: Encoding convolutional activations with deep generative model. In *CVPR*, 2017.
- [15] Z. Qiu, T. Yao, and T. Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In *ICCV*, 2017.
- [16] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel. Self-critical sequence training for image captioning. *arXiv preprint arXiv:1612.00563*, 2016.
- [17] J. B. Roerdink and A. Meijster. The watershed transform: Definitions, algorithms and parallelization strategies. *Fundamenta informaticae*, 2000.
- [18] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015.
- [19] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015.
- [20] H. Wang and C. Schmid. Action recognition with improved trajectories. In *ICCV*, 2013.
- [21] L. Wang, Y. Xiong, Z. Wang, and Y. Qiao. Towards good practices for very deep two-stream convnets. *arXiv preprint arXiv:1507.02159*, 2015.
- [22] T. Yao, T. Mei, and Y. Rui. Highlight detection with pairwise deep ranking for first-person video summarization. In *CVPR*, 2016.
- [23] T. Yao, Y. Pan, Y. Li, and T. Mei. Incorporating copying mechanism in image captioning for learning novel objects. In *CVPR*, 2017.
- [24] T. Yao, Y. Pan, Y. Li, Z. Qiu, and T. Mei. Boosting image captioning with attributes. In *ICCV*, 2017.
- [25] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, D. Lin, and X. Tang. Temporal action detection with structured segment networks. *arXiv preprint arXiv:1704.06228*, 2017.

Activity-Net 2017 Challenge Report

Xiyang Dai
UMD

xdai@umiacs.umd.edu

Bharat Singh
UMD

bharat@umiacs.umd.edu

Larry S. Davis
UMD

lsd@umiacs.umd.edu

In this year’s competition, we use Temporal Context Network (TCN) for precise temporal localization of human activities. To improve performance for the metric used for evaluating proposals (i.e. area under the Average Recall vs. Average Number of Proposals per Video (AR-AN) curve with 100 proposals), we study the influence of two major hyper-parameters: the IOU threshold for NMS and the number of proposals needed from each level of the pyramid. Table 1 shows the impact of NMS IOU threshold on final proposal performance. In contrast to the localization task, we found high IOUs lead to better results on the proposal task. Figure 1 shows the impact of proposals sampled from different anchor lengths (we use percentage of video length to build a temporal pyramid instead of fixed number of frames used in previous works). Table 2 shows our performance on the testing server. After tuning the NMS threshold and proposal anchor length used in TCN, we obtain an AUC score of 61.56 on the testing server, which is ranked third on the ActivityNet 2017 Challenge for the proposal task.

For maximizing recall using a fixed number of proposals, it is important to change the NMS threshold. Since we are allowed to generate 100 proposals (even for 1-2 activities per video), changing NMS is important to improve recall at higher overlap thresholds. Note that if recall is the only evaluation metric, for different numbers of proposals (like 5,10,50,500 etc.) one should pick a different NMS threshold. For example, if recall at only 5 proposals is measured, an NMS threshold of 0.85 would be very bad because most of the proposals would be on the same activity interval.

Please check out the full paper¹ for further details on TCN.

NMS	0.3	0.45	0.6	0.75	0.8	0.85	0.9
AUC	42.52	47.14	55.44	57.94	58.22	58.54	57.98

Table 1: Impact of NMS IOU threshold on the proposal performance on the ActivityNet validation set

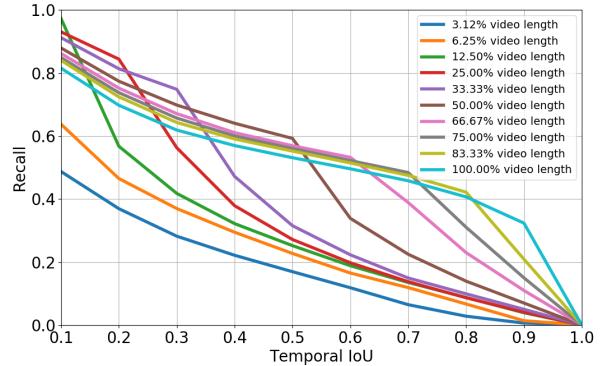


Figure 1: Impact of different anchor lengths on proposal performance on ActivityNet validation set

NMS threshold	AUC on Validation	AUC on Test
baseline	47.14	49.03
fine-tune NMS	58.54	59.89
final	59.58	61.56

Table 2: Our final proposal performance on ActivityNet 2017 Challenge

¹ <https://arxiv.org/pdf/1708.02349.pdf>

Temporal Convolution Based Action Proposal: Submission to ActivityNet 2017

Tianwei Lin¹, Xu Zhao^{1*}, Zheng Shou²

¹Computer Vision Laboratory, Shanghai Jiao Tong University, China. ²Columbia University, USA

{wzmsltw, zhaoxu}@sjtu.edu.cn, zs2262@columbia.edu

Abstract

In this notebook paper, we describe our approach in the submission to the temporal action proposal (task 3) and temporal action localization (task 4) of ActivityNet Challenge hosted at CVPR 2017. Since the accuracy in action classification task is already very high (nearly 90% in ActivityNet dataset), we believe that the main bottleneck for temporal action localization is the quality of action proposals. Therefore, we mainly focus on the temporal action proposal task and propose a new proposal model based on temporal convolutional network. Our approach achieves the state-of-the-art performances on both temporal action proposal task and temporal action localization task.

1. Introduction

Action recognition and temporal action localization are both important branches of video content analysis. The temporal action localization or detection task aims to detect action instances in untrimmed video, containing categories and temporal boundaries of action instances.

Temporal action localization task can be divided into two main parts. (1) Temporal action proposal, which means we need generate some temporal boundaries of action instances without classifying their categories. (2) Action recognition (or we can say action classification), in this part we need to decide the categories of temporal action proposals. Most of previous works [11, 15] address these two parts separately. There are also works [3, 1] focusing on temporal action proposal. For action recognition, there are already many algorithms [12, 4] with great performance. However, the localization accuracy (mean average precision) is still very low in multiple benchmarks such as THUMOS'14 [6] and ActivityNet [2], comparing with the situation in object localization. We think the main constraint on accuracy of temporal action localization is the quality of action proposals. Therefore, we mainly focus on the temporal action proposal task in this challenge and our high quality proposals

also lead to state-of-the-art performance in temporal action localization task.

2. Our Approach

The framework of our approach is shown in Fig 1. In this section, we introduce each part of the framework, which consists of feature extraction, temporal action proposal and temporal action localization.

2.1. Feature Extraction

The first step of our framework is feature extraction. We extract two-stream features in a similar way described in [5]. We adopt two-stream network [14] which is pre-trained on ActivityNet v1.3 training set. First we segment video into 16-frames snippets without overlap. In each snippet, we use spatial network to extract appearance feature with central frame, and we use the output of “Flatten-673” layer in ResNet network as feature. For motion feature, we compute optical flows using 6 consecutive frames around the center frame of a snippet, then these optical flows are used for extracting motion feature with temporal network, where the output of “global-pool” layer in BN-Inception network is used as feature. Then, we concatenate appearance and motion feature to form the snippet-level features, which are 3072-dimensional vectors. So after feature extraction, we can transfer a video into a sequence of snippet-level feature vectors. Finally, we resize the feature sequence to new length 256 by linear interpolation.

Since we only use two-stream network trained on ActivityNet v1.3 training set to extract features, there is no external data used in our approach.

2.2. Temporal Action Proposal

Prop-SSAD. In our previous work [7]¹, we design a model called **Single Shot Action Detector (SSAD)** network which simultaneously conducts temporal action proposal and recognition. A core idea of SSAD is applying anchor mechanism to temporal action localization task based on temporal convolutional layers, which is similar

*Corresponding author.

¹This paper can be found at: <https://wzmsltw.github.io/>

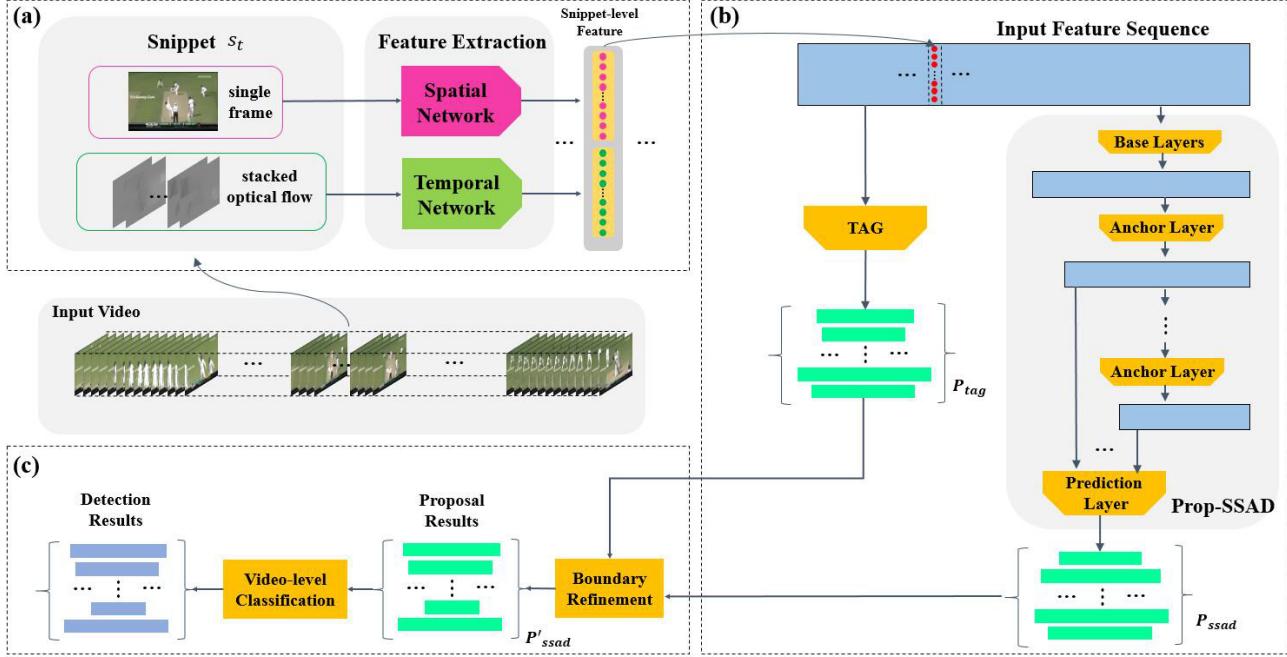


Figure 1: The framework of our approach. (a) Two-stream networks are used to extract snippet-level features. (b) Prop-SSAD model and TAG method are used for proposal generation separately. (c) Proposals generated by TAG are used for refining the boundaries of proposals generated by Prop-SSAD model. We use video-level action classification result as the category of temporal action proposals to get temporal action localization result.

with YOLO [9] and SSD [8] network for object localization task. In detail, we associate multiple temporal anchor instances with multi-scale temporal feature maps, then use temporal convolutional layers to predict information of anchor instances, including action categories, overlap score and location offsets. So SSAD can directly detect temporal action instances using feature sequence of untrimmed video.

In this challenge, we use SSAD network to make temporal action proposal without action recognition and we call it **Prop-SSAD**. The main differences of network configuration between Prop-SSAD and SSAD are listed below.

- Type of input features. In SSAD, we use two-stream network and C3D network to extract feature of video; in Prop-SSAD, only two-stream networks are used.
- Number of anchor layers. In SSAD, we only associate temporal anchors with 3 temporal feature maps using anchor layers with length 4, 8 and 16; in Prop-SSAD, we associate temporal anchors with 7 temporal feature maps with length 1, 2, 4, 8, 16, 32, 64.
- Loss function. In SSAD, we use classification, overlap and location loss jointly to train network; in Prop-SSAD, only overlap loss is used.

TAG [15]. We also implement **Temporal Actionness Grouping (TAG)** method to generate temporal action proposals, which is proposed in [15]. Since the code of TAG is not released yet, we implement TAG by ourselves. First we train a multi-layer perceptron (MLP) model with one hidden layer to predict the actionness score for each snippet, then we use grouping method described in [15] with multiple threshold to generate temporal action proposals. Proposals generated by TAG are used for refining the proposals' boundaries generated by Prop-SSAD.

Boundaries Refinement. Given feature vector sequence of a video, we can get temporal action proposals set P_{ssad} using Prop-SSAD and temporal action proposals set P_{tag} using TAG. For each proposal p_t in P_{tag} , we calculate its IoU with all proposals in P_{ssad} . If the maximum IoU is higher than threshold 0.75, we replace the boundaries of corresponding proposal p_s in P_{ssad} with boundaries of p_t . After refinement procedure, we get refined proposals set P'_{ssad} , which is the final proposal results.

2.3. Temporal Action Localization

Since most videos in ActivityNet dataset only contain one action category, we use video-level action classification result as the category of temporal action proposals to get temporal action localization result.

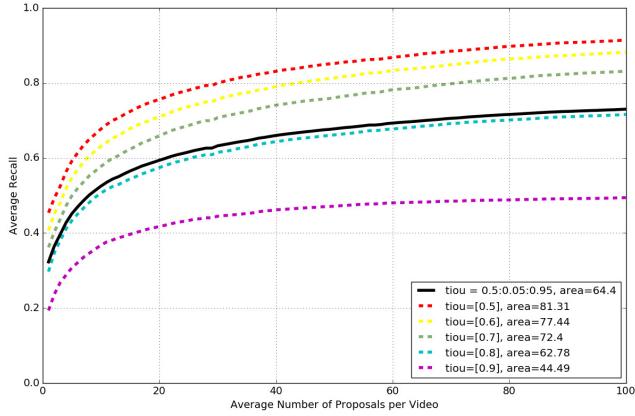


Figure 2: AR-AN curve of our proposal results in validation set. The area under black curve is the AR-AN score.

Table 1: Proposal Results on validation set of ActivityNet.

Method	AR@10	AR@100	AR-AN
Uniform Random (baseline)	29.02	55.71	44.88
Prop-SSAD	50.44	69.54	61.52
Refined Prop-SSAD	52.50	73.01	64.40

3. Experimental Results

3.1. Evaluation Metrics

Localization. In temporal action localization task, mean Average Precision (mAP) is used as the metric for evaluating result, which is similar with metrics used in object localization task. In detail, the official metric used in this task is the average mAP computed with tIoU thresholds between 0.5 and 0.95 with the step size of 0.05.

Proposal. In temporal action proposal task, the area under the Average Recall vs. Average Number of Proposals per Video (**AR-AN**) curve is used as the evaluation metric, where **AR** is defined as the mean of all recall values using tIoU thresholds between 0.5 and 0.95 with a step size of 0.05. In this notebook, we call **AR** with a certain number of **AN** as **AR@AN**. For example, AR@100 means average recall with 100 proposals.

3.2. Temporal Action Proposal

The proposal performance on validation set of our approach are shown in Table 1 and Figure 2. Our approach significantly outperform the baseline method and refined Prop-SSAD has better performance than Prop-SSAD. The boundaries refinement mainly improve the average recall with high tIoU.

Table 2: Action localization results on validation set. Results are evaluated by mAP with different IoU thresholds α and average mAP of IoU thresholds from 0.5 to 0.95. Ours@n means first n proposals used for localization.

mAP	0.5	0.75	0.95	Average mAP
Wang et al. [10]	42.28	3.76	0.05	14.85
Shou et al. [10]	43.83	25.88	0.21	22.77
Xiong et. al. [15]	39.12	23.48	5.49	23.98
Ours@1	42.14	27.17	6.54	27.00
Ours@5	46.56	30.94	7.53	30.49
Ours@10	47.84	31.90	7.76	31.41
Ours@25	48.56	32.53	7.83	31.93
Ours@100	48.99	32.91	7.87	32.26

Table 3: Action localization results on testing set. Only average mAP is provided in evaluation server, which is calculated with IoU thresholds from 0.5 to 0.95.

Method	Average mAP
Wang et. al. [13]	14.62
Xiong et. al. [15]	26.05
Zhao et. al. [16]	28.28
Ours result	33.40

3.3. Temporal Action Localization

In the temporal action localization task, we directly use proposals submitted in temporal action proposal task. For action categories, we use the video-level classification results of [13], which has 87.7% Top-1 classification accuracy and obtains 2nd place in untrimmed video classification task of ActivityNet Challenge 2016.

Evaluation results in validation set are shown in Table 2. These results suggest that localization mAP mainly depends on first several proposals. Therefore, we think AR-AN may not be the best evaluation metric for temporal action proposal task. AR with small proposals amount should has higher weight in evaluation metric.

Evaluation results in testing set are shown in Table 3. Our approach significantly outperform other state-of-the-art approaches. We think the main contributor is our high quality temporal action proposals.

4. Conclusion

In this challenge, we mainly focus on the temporal action proposal task and obtains the salient performance in both temporal action proposal and temporal action localization task. Our results suggested that anchor mechanisms and temporal convolution can work well in temporal action proposal task. In the future, we will improve our framework such as training the whole networks end-to-end.

References

- [1] F. Caba Heilbron, J. Carlos Niebles, and B. Ghanem. Fast temporal activity proposals for efficient detection of human actions in untrimmed videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1914–1923, 2016.
- [2] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
- [3] V. Escorcia, F. C. Heilbron, J. C. Niebles, and B. Ghanem. Daps: Deep action proposals for action understanding. In *European Conference on Computer Vision*, pages 768–784. Springer, 2016.
- [4] C. Feichtenhofer, A. Pinz, and A. Zisserman. Convolutional two-stream network fusion for video action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1933–1941, 2016.
- [5] J. Gao, Z. Yang, and R. Nevatia. Cascaded boundary regression for temporal action detection. *arXiv preprint arXiv:1705.01180*, 2017.
- [6] Y. G. Jiang, J. Liu, A. R. Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. Thumos challenge: Action recognition with a large number of classes. In *ECCV Workshop*, 2014.
- [7] T. Lin, X. Zhao, and Z. Shou. Single shot temporal action detection. *25nd ACM international conference on Multimedia*, 2017.
- [8] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European Conference on Computer Vision*, pages 21–37. Springer, 2016.
- [9] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788, 2016.
- [10] Z. Shou, J. Chan, A. Zareian, K. Miyazawa, and S.-F. Chang. Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. *arXiv preprint arXiv:1703.01515*, 2017.
- [11] Z. Shou, D. Wang, and S.-F. Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1049–1058, 2016.
- [12] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: towards good practices for deep action recognition. In *European Conference on Computer Vision*, pages 20–36. Springer, 2016.
- [13] R. Wang and D. Tao. Uts at activitynet 2016. *ActivityNet Large Scale Activity Recognition Challenge*, 2016:8, 2016.
- [14] Y. Xiong, L. Wang, Z. Wang, B. Zhang, H. Song, W. Li, D. Lin, Y. Qiao, L. V. Gool, and X. Tang. Cuhk & ethz & siat submission to activitynet challenge 2016. *arXiv preprint arXiv:1608.00797*, 2016.
- [15] Y. Xiong, Y. Zhao, L. Wang, D. Lin, and X. Tang. A pursuit of temporal accuracy in general activity detection. *arXiv preprint arXiv:1703.02716*, 2017.
- [16] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, D. Lin, and X. Tang. Temporal action detection with structured segment networks. *arXiv preprint arXiv:1704.06228*, 2017.

R-C3D: Region Convolutional 3D Network for Temporal Activity Detection

Huijuan Xu

Abir Das
Boston University
Boston, MA

Kate Saenko

{hxu, dasabir, saenko}@bu.edu

Abstract

We address the problem of activity detection in continuous, untrimmed video streams. This is a difficult task that requires extracting meaningful spatio-temporal features to capture activities, accurately localizing the start and end times of each activity, and also dealing with very large data volumes. We introduce a new model, *Region Convolutional 3D Network (R-C3D)*, which encodes the video streams using a three-dimensional fully convolutional network, then generates candidate temporal regions containing activities, and finally classifies selected regions into specific activities. Computation is saved due to the sharing of convolutional features between the proposal and the classification pipelines. The entire model is trained end-to-end with jointly optimized localization and classification losses. *R-C3D* is faster than existing methods (569 frames per second on a single Titan X Maxwell GPU) and achieves state-of-the-art results on THUMOS'14 (10% absolute improvement). We further demonstrate that our model is a general activity detection framework that does not rely on assumptions about particular dataset properties by evaluating our approach on ActivityNet and Charades.

1. Introduction

Activity detection in continuous video is a challenging problem that requires not only recognizing, but also precisely localizing activities in time. Existing state-of-the-art approaches address this task as *detection by classification*, *i.e.* classifying temporal segments generated in the form of sliding windows [13, 20, 24, 37] or via an external “proposal” generation mechanism [10, 35]. These approaches suffer from one or more of the following major drawbacks: they do not learn deep representations in an end-to-end fashion, but rather use hand-crafted features [33, 34], or deep features like VGG [28], ResNet [8], C3D [32] etc., learned separately on image/video classification tasks. Such off-the-shelf representations may not be optimal for localizing activities in diverse video domains, resulting in inferior performance. Furthermore, current methods’ dependence

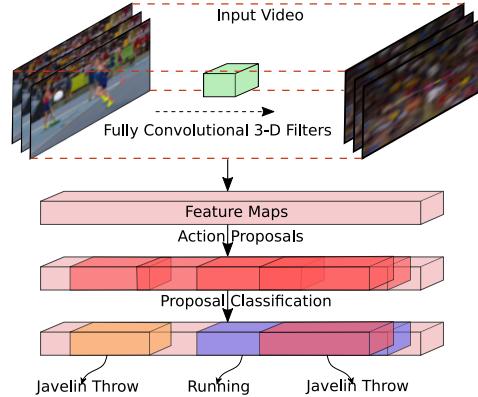


Figure 1. We propose a fast end-to-end *Region Convolutional 3D Network (R-C3D)* for activity detection in continuous video streams. The network encodes the frame buffer with fully-convolutional 3D filters, proposes activity segments, then classifies and refines them based on pooled features within their boundaries. Our model improves both speed and accuracy compared to existing methods.

on external proposal generation or exhaustive sliding windows leads to poor computational efficiency. Finally, the sliding-window models cannot easily predict flexible activity boundaries.

In this paper, we propose an activity detection model that addresses all of the above issues. Our *Region Convolutional 3D Network (R-C3D)* is end-to-end trainable and learns task-dependent convolutional features by jointly optimizing proposal generation and activity classification. Inspired by the Faster R-CNN [21] object detection approach, we compute fully-convolutional 3D ConvNet features and propose temporal regions likely to contain activities, then pool features within these 3D regions to predict activity classes (Figure 1). The proposal generation stage filters out many background segments and results in superior computational efficiency compared to sliding window models. Furthermore, proposals are predicted with respect to predefined anchor segments and can be of arbitrary length, allowing detection of flexible activity boundaries.

Convolutional Neural Network (CNN) features learned end-to-end have been successfully used for activity recog-

nition [14, 27], particularly in 3D ConvNets (C3D [32]), which learn to capture spatio-temporal features. However, unlike the traditional usage of 3D ConvNets [32] where the input is short 16-frame video chunks, our method applies full convolution along the temporal dimension to encode as many frames as the GPU memory allows. Thus, rich spatio-temporal features are automatically learned from longer videos. These feature maps are shared between the activity proposal and classification subnets to save computation time and jointly optimize features for both tasks.

Alternative activity detection approaches [4, 17, 18, 29, 39] use a recurrent neural network (RNN) to encode a sequence of frame or video chunk features (*e.g.* VGG [28], C3D [32]) and predict the activity label at each time step. However, these RNN methods can only model temporal features at a fixed granularity (*e.g.* per-frame CNN features or 16-frame C3D features). In order to use the same classification network to classify variable length proposals into specific activities, we extend 2D region of interest (RoI) pooling to 3D which extracts a fixed-length feature representation for these proposals. Thus, our model can utilize video features at any temporal granularity. Furthermore, some RNN-based detectors rely on direct regression to predict the temporal boundaries for each activity. As shown by results in object detection [7, 31] and semantic segmentation [2], using a regression-only framework to predict object boundaries does not work well in practice compared to “proposal based detection”.

We perform extensive comparisons of our approach to state-of-the-art activity detection methods using three publicly available benchmark datasets - THUMOS’14 [12] sports activities, ActivityNet [9] human activities on Youtube, and Charades [26] in-the-wild daily activities. We achieve new state-of-the-art results on THUMOS’14 and Charades, and improved results on ActivityNet when using only C3D features. Our code will be made publicly available to support further research progress.

To summarize, the main contributions of our paper are:

- an end-to-end activity detection model with combined activity proposal and classification stages that can detect arbitrary length activities;
- fast detection speeds (5x faster than current methods) achieved by sharing fully-convolutional C3D features between the proposal generation and classification parts of the network;
- extensive evaluations on three diverse activity detection datasets that demonstrate the general applicability of our model.

2. Related Work

Activity Detection There is a long history of activity recognition, or classifying trimmed video clips into fixed set of categories [11, 15, 19, 27, 33, 42]. Activity *detection*

also needs to predict the start and end times of the activities within untrimmed and long videos. Existing activity detection approaches are dominated by models that use sliding windows to generate segments and subsequently classify them with activity classifiers trained on multiple features [13, 20, 24, 37]. Most of these methods have stage-wise pipelines which are not trained end-to-end. Moreover, the use of exhaustive sliding windows is computationally inefficient and constrains the boundary of the detected activities to some extent.

Recently, some approaches have bypassed the need for exhaustive sliding window search to detect activities with arbitrary lengths. [4, 17, 18, 29, 39] achieve this by modeling the temporal evolution of activities using RNNs or Long Short Term Memory (LSTM) networks and predicting an activity label at each time step. The deep action proposal model in [4] uses LSTM to encode C3D features of every 16-frame video chunk, and directly regresses and classifies activity segments without the extra proposal generation stage. Compared to this work, we avoid recurrent layers, encoding a large video buffer with a fully-convolutional 3D ConvNet, and use 3D RoI pooling to allow feature extraction at arbitrary proposal granularity, achieving significantly higher accuracy and speed. The method in [41] tries to capture motion features at multiple resolutions by proposing a Pyramid of Score Distribution Features for activity detection, however their model is not end-to-end trainable and relies on handcrafted features.

Aside from supervised activity detection, a recent work [36] has addressed weakly supervised activity localization from training data labeled only with video level class labels by learning attention weights on shot based or uniformly sampled proposals. The framework proposed in [22] explores the uses of a language model and an activity length model in a detection pipeline. Spatio-temporal activity localization [38, 40] have also been explored to some extent. We only focus on supervised temporal activity localization in this work.

Object Detection Activity detection in untrimmed videos is intricately related to object detection in images. The inspiration for our work, Faster R-CNN [21], extends R-CNN [7] and Fast R-CNN [6] object detection approaches, incorporating RoI pooling and a region proposal network. Compared to recent object detection models *e.g.*, SSD [16] and R-FCN[3], Faster R-CNN is a general and robust object detection framework that has been deployed on different datasets with little data augmentation effort. Like Faster R-CNN, our R-C3D model is also designed with the goal of easy deployment on varied activity detection datasets. It avoids making certain assumptions based on unique characteristics of a dataset, such as the UPC model for ActivityNet [18] which assumes that each video contains a single activity class. We show the effectiveness of our

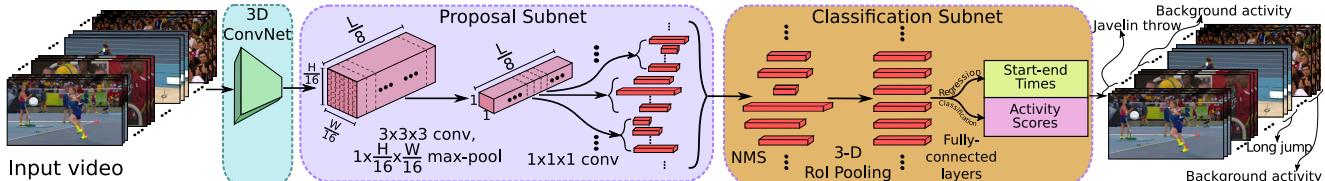


Figure 2. R-C3D model architecture. The 3D ConvNet takes raw video frames as input and computes convolutional features. These are input to the Proposal Subnet that proposes candidate activities of variable length along with confidence scores. The Classification Subnet filters the proposals, pools fixed size features and then predicts activity labels along with refined segment boundaries.

model on three different types of activity detection datasets, the most extensive evaluation to our knowledge.

3. Approach

We propose a *Region Convolutional 3D Network (R-C3D)*, a novel convolutional neural network for activity detection in continuous video streams. The network, illustrated in Figure 2, consists of three components: a shared 3D ConvNet feature extractor [32], a temporal proposal stage, and an activity classification and refinement stage. To enable efficient computation and end-to-end training, the proposal and classification sub-networks share the same C3D feature maps. The proposal subnet predicts variable length temporal segments that potentially contain activities, while the classification subnet classifies these proposals into specific activity categories or background, and further refines the proposal segment boundaries. A key innovation is to extend the 2D ROI pooling in Faster R-CNN to 3D ROI pooling which allows our model to extract features at various resolutions for variable length proposals. Next, we describe the shared video feature hierarchies in Sec. 3.1, the temporal proposal subnet in Sec. 3.2 and the classification subnet in Sec. 3.3. Sections 3.4 and 3.5 detail the optimization strategy during training and testing respectively.

3.1. 3D Convolutional Feature Hierarchies

We use a 3D ConvNet to extract rich spatio-temporal feature hierarchies from a given input video buffer. It has been shown that both spatial and temporal features are important for representing videos, and a 3D ConvNet encodes rich spatial and temporal features in a hierarchical manner. The input to our model is a sequence of RGB video frames with dimension $\mathbb{R}^{3 \times L \times H \times W}$. The architecture of the 3D ConvNet is taken from the C3D architecture proposed in [32]. However, unlike [32], the input to our model is of variable length. We adopt the convolutional layers (from conv1a to conv5b) of C3D, so a feature map $C_{conv5b} \in \mathbb{R}^{512 \times \frac{L}{8} \times \frac{H}{16} \times \frac{W}{16}}$ (512 is the channel dimension of the layer conv5b) is produced as the output of this sub-network. We use C_{conv5b} activations as the shared input to the proposal and classification subnets. The height (H) and width (W) of the frames are taken as 112 each follow-

ing [32]. The number of frames L can be arbitrary and is only limited by memory.

3.2. Temporal Proposal Subnet

To allow the model to predict variable length proposals, we incorporate anchor segments into the temporal proposal sub-network. The subnet predicts potential proposal segments with respect to anchor segments and a binary label indicating whether the predicted proposal contains an activity or not. The anchor segments are pre-defined multiscale windows centered at $(L/8)$ uniformly distributed temporal locations. Each temporal location specifies K anchor segments, each at a different fixed scale. Thus, the total number of anchor segments is $(L/8) * K$. The same set of K anchor segments exists in different temporal locations, which ensures that the proposal prediction is temporally invariant. The anchors serve as reference activity segments for proposals at each temporal location, where the maximum number of scales K is dataset dependent.

To obtain features at each temporal location for predicting proposals with respect to these anchor segments, we first add a 3D convolutional filter with kernel size $3 \times 3 \times 3$ on top of C_{conv5b} to extend the temporal receptive field for the temporal proposal subnet. Then, we downsample the spatial dimensions (from $\frac{H}{16} \times \frac{W}{16}$ to 1×1) to produce a *temporal* only feature map $C_{tpn} \in \mathbb{R}^{512 \times \frac{L}{8} \times 1 \times 1}$ by applying a 3D max-pooling filter with kernel size $1 \times \frac{H}{16} \times \frac{W}{16}$. The 512-dimensional feature vector at each temporal location in C_{tpn} is used to predict a relative offset $\{\delta c_i, \delta l_i\}, i \in \{1, \dots, K\}$ to the center location and the length of each anchor segment $\{c_i, l_i\}, i \in \{1, \dots, K\}$. It also predicts the binary scores for each proposal being an activity or background. The proposal offsets and scores are predicted by adding two $1 \times 1 \times 1$ convolutional layers on top of C_{tpn} .

Training: For training, we need to assign positive/negative labels to the anchor segments. Following the standard practice in object detection [21], we choose a positive label if the anchor segment 1) overlaps with some ground-truth activity with Intersection-over-Union (IoU) higher than 0.7, or 2) has the highest IoU overlap with some ground-truth activity. If the anchor has IoU overlap lower than 0.3 with all ground-truth activities, then it is given a negative label. All others are held out from training. For proposal regression,

the anchor segment is transformed with respect to a nearby ground truth activity segment using the coordinate transformations described in Sec. 3.4. We sample balanced batches with a positive/negative ratio of 1:1.

3.3. Activity Classification Subnet

The activity classification stage has three main functions: 1) selecting proposal segments from the previous stage, 2) three-dimensional region of interest (3D RoI) pooling to extract fixed-size features for selected proposals, and 3) activity classification and boundary regression for the selected proposals based on the pooled features.

The proposal subnet outputs a set of candidate proposal segments with associated scores. Some activity proposals highly overlap with each other and some have a low proposal score indicating low confidence. Following the standard practice in object detection [5, 21] and activity detection [24, 39], we employ a greedy Non-Maximum Suppression (NMS) strategy to eliminate highly overlapping and low confidence proposals. The NMS threshold is set as 0.7.

The selected proposals can be of arbitrary length. However we need to extract fixed-size features for each of them in order to use fully connected layers for further activity classification and regression. We design a 3D RoI pooling layer to extract the fixed-size volume features for each variable-length proposal from the shared convolutional features $C_{conv5b} \in \mathbb{R}^{512 \times (L/8) \times 7 \times 7}$ (shared with the temporal proposal subnet). Specifically, in 3D RoI pooling, an input feature volume of size, say, $l \times h \times w$ is divided into $l_s \times h_s \times w_s$ sub-volumes each with approximate size $\frac{l}{l_s} \times \frac{h}{h_s} \times \frac{w}{w_s}$, and then max pooling is performed inside each sub-volume. In our case, suppose a proposal has the feature volume of $l_p \times 7 \times 7$ in C_{conv5b} , then this feature volume will be divided into $1 \times 4 \times 4$ grids and max pooled inside each grid. Thus, proposals of arbitrary lengths give rise to output volume features of the same size $512 \times 1 \times 4 \times 4$.

The output of the 3D RoI pooling for selected proposals is fed to a series of two fully connected layers. Finally the proposals are classified to activity categories by a classification layer and the refined start-end times for these proposals are given by a regression layer. The classification and regression layers are also two separate fully connected layers and for both of them the input comes from the aforementioned fully connected layers (after the 3D RoI pooling layer).

Training: We need to assign an activity label to each proposal predicted by the proposal subnet for training the classifier subnet. An activity label is assigned if the proposal has the highest IoU overlap with a ground-truth activity, and at the same time, the IoU overlap is greater than 0.5. A background label (no activity) is assigned to proposals with IoU overlap lower than 0.5 with all ground-truth activities. Training batches are chosen with positive/negative

ratio of 1:3.

3.4. Optimization

We train the network by optimizing both the classification and regression tasks jointly for the two subnets. The softmax loss function is used for classification, and smooth L1 loss function [6] is used for regression. Specifically, the objective function is given by:

$$Loss = \frac{1}{N_{cls}} \sum_i L_{cls}(a_i, a_i^*) + \lambda \frac{1}{N_{reg}} \sum_i a_i^* L_{reg}(t_i, t_i^*) \quad (1)$$

where N_{cls} and N_{reg} stand for batch size and the number of anchor/proposal segments, λ is the loss trade-off parameter and is set to a value 1. i is the anchor/proposal segments index in a batch, a_i is the predicted probability of the proposal or activities, a_i^* is the ground truth, $t_i = \{\delta c_i, \delta l_i\}$ represents predicted relative offset to anchor segments or proposals. $t_i^* = \{\delta c_i^*, \delta l_i^*\}$ represents the coordinate transformation of ground truth segments to anchor segments or proposals. The coordinate transformations are computed as follows:

$$\begin{cases} \delta c_i = (c_i^* - c_i)/l_i \\ \delta l_i = \log(l_i^*/l_i) \end{cases} \quad (2)$$

where c_i and l_i are the center location and the length of anchor segments or proposals while c_i^* and l_i^* denote the same for the ground truth activity segments.

In our R-C3D model, the above loss function is applied for both the temporal proposal subnet and the activity classification subnet. In the proposal subnet, the binary classification loss L_{cls} predicts whether the proposal contains an activity or not, and the regression loss L_{reg} optimizes the relative displacement between proposals and ground truths. In the proposal subnet the losses are activity class agnostic. For the activity classification subnet, the multiclass classification loss L_{cls} predicts the specific activity class for the proposal, and the number of classes are the number of activities plus one for the background. The regression loss L_{reg} optimizes the relative displacement between activities and ground truths. All four losses for the two subnets are optimized jointly.

3.5. Prediction

Activity prediction in R-C3D consists of two steps. First, the proposal subnet generates candidate proposals and predicts the start-end time offsets as well as proposal score for each. Then the proposals are refined via NMS with threshold value 0.7. After NMS, the selected proposals are fed to the classification network to be classified into specific activity classes, and the activity boundaries of the predicted proposals are further refined by the regression layer. The localization prediction in both proposal subnet and classification subnet is in the form of relative displacement of center point and length of segments. In order to get the start time

Table 1. Activity detection results on THUMOS’14 (in percentage). mAP at different IoU thresholds α are reported. The top three performers on the THUMOS’14 challenge leaderboard and other results reported in existing papers are shown.

	α				
	0.1	0.2	0.3	0.4	0.5
Karaman et al. [13]	4.6	3.4	2.1	1.4	0.9
Wang et al. [37]	18.2	17.0	14.0	11.7	8.3
Oneata et al. [20]	36.6	33.6	27.0	20.8	14.4
Heilbron et al. [10]	-	-	-	-	13.5
Escorcia et al. [4]	-	-	-	-	13.9
Richard et al. [22]	39.7	35.7	30.0	23.2	15.2
Yeung et al. [39]	48.9	44.0	36.0	26.4	17.1
Yuan et al. [41]	51.4	42.6	33.6	26.1	18.8
Shou et al. [24]	47.7	43.5	36.3	28.7	19.0
Shou et al. [23]	-	-	40.1	29.4	23.3
R-C3D (our one-way buffer)	51.6	49.2	42.8	33.4	27.0
R-C3D (our two-way buffer)	54.5	51.5	44.8	35.6	28.9

and end time of the predicted proposals or activities, inverse coordinate transformation to Equation 2 is performed.

Our model accepts variable length input videos. However, to take advantage of the vectorized implementation in fast deep learning libraries, we pad the last few frames of short videos with blank frames, and break long videos into buffers (limited by memory only). The predicted activities are post-precessed by NMS at a lower threshold (0.1 lower than the mAP evaluation threshold) to get the final activity predictions.

4. Experiments

We evaluate R-C3D on three large-scale activity detection datasets - THUMOS’14 [12], Charades[26] and ActivityNet [9]. Sections 4.1, 4.2, 4.3 provide the experimental details and evaluation results on these three datasets. Results are shown in terms of mean Average Precision - mAP@ α where α denotes different Intersection over Union (IoU) thresholds, as is the common practice in the literature. Section 4.4 provides the detection speed comparison with state-of-the-art methods.

4.1. Experiments on THUMOS’14

The THUMOS’14 activity detection dataset contains over 24 hours of video from 20 different sport activities. The training set contains 2765 trimmed videos while the validation set and the test set contain 200 and 213 untrimmed videos respectively. This dataset is particularly challenging as it consists of very long videos (up to a few hundreds of seconds) with multiple activity instances of very small duration (up to few tens of seconds). Most videos contain multiple activity instances of the same activity class. In addition, some videos contain activity segments from different classes.

Experimental Setup: We divide 200 untrimmed videos from the validation set into 180 training and 20 held out

Table 2. Per-class AP at IoU threshold $\alpha = 0.5$ on THUMOS’14 (in percentage).

	[20]	[39]	[24]	R-C3D (ours)
Baseball Pitch	8.6	14.6	14.9	26.1
Basketball Dunk	1.0	6.3	20.1	54.0
Billiards	2.6	9.4	7.6	8.3
Clean and Jerk	13.3	42.8	24.8	27.9
Cliff Diving	17.7	15.6	27.5	49.2
Cricket Bowling	9.5	10.8	15.7	30.6
Cricket Shot	2.6	3.5	13.8	10.9
Diving	4.6	10.8	17.6	26.2
Frisbee Catch	1.2	10.4	15.3	20.1
Golf Swing	22.6	13.8	18.2	16.1
Hammer Throw	34.7	28.9	19.1	43.2
High Jump	17.6	33.3	20.0	30.9
Javelin Throw	22.0	20.4	18.2	47.0
Long Jump	47.6	39.0	34.8	57.4
Pole Vault	19.6	16.3	32.1	42.7
Shotput	11.9	16.6	12.1	19.4
Soccer Penalty	8.7	8.3	19.2	15.8
Tennis Swing	3.0	5.6	19.3	16.6
Throw Discus	36.2	29.5	24.4	29.2
Volleyball Spiking	1.4	5.2	4.6	5.6
mAP@0.5	14.4	17.1	19.0	28.9

videos to get the best hyperparameter setting. All 200 videos are used as the training set and the final results are reported on 213 videos in the test set. Since the GPU memory is limited, we first create a buffer of 768 frames at 25 frames per second (fps) which means approximately 30 seconds of video. Our choice is motivated by the fact that 99.5% of all activity segments in the validation set (used here as the training set) are less than 30 seconds long. These buffers of frames act as inputs to R-C3D . We can create the buffer by sliding from the beginning of the video to the end, denoted as the ‘one-way buffer’. An additional pass from the end of the video to the beginning can be used to increase the amount of training data as a data augmentation strategy, denoted as the ‘two-way buffer’. We initialize the 3D ConvNet part of our model with C3D weights trained on Sports-1M and finetuned on UCF101 released by the author in [32]. We allow all the layers of R-C3D to be trained on THUMOS’14 with a fixed learning rate of 0.0001.

The number of anchor segments K chosen for this dataset is 10 with specific scale values of [2, 4, 5, 6, 8, 9, 10, 12, 14, 16]. The values are chosen according to the distribution of the activity durations in the training set. At 25 fps and temporal pooling factor of 8 (C_{tpn} downsamples the input by 8 temporally), the anchor segments correspond to segments of duration between 0.64 and 5.12 seconds¹. Note that, the predicted proposals or activities are relative to the anchor segments but not limited to the anchor segment boundaries, enabling our model to detect variable-length activities.

¹ $2 * 8/25 = 0.64$ and $16 * 8/25 = 5.12$

Results: As a sanity check, we first evaluate the performance of the temporal proposal subnet (ref Section 3.2). A predicted proposal is marked as correct if it has IoU with a ground truth activity of more than 0.7, otherwise it is considered incorrect. With this binary setting, precision and recall values for the temporal proposal subnet are relatively high - 85% and 83% respectively.

In Table 1, we present a comparative evaluation of the activity detection performance of our end-to-end model with existing state-of-the-art approaches in terms of mAP at IoU thresholds 0.1-0.5 (denoted as α). For both the one-way buffer setting and the two-way buffer setting we achieve new state-of-the-art for all five α values. In the one-way setting, mAP@0.5 is 27.0% which is an 8% absolute improvement from the state-of-the-art. The two-way buffer setting further increases the mAP values at all the IoU thresholds with mAP@0.5 reaching as far as 28.9%. Our model comprehensively outperforms the current state-of-the-art by a large margin (28.9% compared to 19.0% as reported in [24]).

The Average Precision (AP) for each class in THUMOS’14 at IoU threshold 0.5 for the two-way buffer setting is shown in Table 2. For per-class AP, our model outperforms the other three baselines in most classes and it shows significant improvement (by more than 20% absolute AP) for activities *e.g.*, Basketball Dunk, Cliff Diving, and Javelin Throw. For some of the activities, our method is only second to the best performing ones by a very small margin (*e.g.*, Billiards or Cricket Shot). Figure 3(a) shows some representative qualitative results from two videos in this dataset.

4.2. Experiments on ActivityNet

The ActivityNet [9] dataset consists only of untrimmed videos and is released in three versions. We use the latest release (ActivityNet 1.3) which has 10024, 4926 and 5044 videos containing 200 different types of activities in the train, validation and test sets respectively. Except for only a few, most videos contain activity instances of a single class covering a great deal of the video. Compared to THUMOS’14, this is a large-scale dataset both in terms of the number of activities involved and the amount of video. Researchers have taken part in the ActivityNet challenge [1] held on this dataset. The performance of the participating teams is evaluated on test videos for which the ground truth annotations are not public. In addition to evaluating on the validation set, we show our performance on the test set after evaluating it on the challenge server.

Experimental Setup: Similar to THUMOS’14, we keep the length of the input buffer to be 768 but, as the videos are long, we sample frames at 3 fps to fit it in the GPU memory. This makes the duration of the buffer approximately 256 seconds which covers over 99.99% activities in the training

Table 3. Detection results on ActivityNet in terms of mAP@0.5 (in percentage). The top half of the table shows performance from methods using additional handcrafted features while the bottom half shows approaches using deep features only (including ours). Results for [29] are taken from [1]

	train data	validation	test
G. Singh <i>et. al.</i> [30]	train	34.5	36.4
B. Singh <i>et. al.</i> [29]	train+val	-	28.8
UPC [18]	train	22.5	22.3
R-C3D (ours)	train	26.8	26.8
R-C3D (ours)	train+val	-	28.4

split. The considerably long activity durations in ActivityNet prompt us to set the number of anchor segments K to be as high as 20. Specifically, we chose the following scales - [1, 2, 3, 4, 5, 6, 7, 8, 10, 12, 14, 16, 20, 24, 28, 32, 40, 48, 56, 64]. Thus the shortest anchor segment is of duration 2.7 seconds and the longest one is of duration 170 seconds, which covers 95.6% of all activities in the training set.

Considering the vast domain difference of the activities between Sports-1M and ActivityNet, we finetune the Sports-1M pretrained 3D ConvNet model [32] with the training videos of ActivityNet at 3 fps on the activity classification task. We initialize the 3D ConvNet part of our model with these finetuned weights. ActivityNet being a large scale dataset, the training takes more epochs. As a speed-efficiency trade-off, we freeze the first two convolutional layers in our model during training. The learning rate is kept fixed at 0.0001 for the first 10 epochs and then it is decreased to 0.00001 for 5 further epochs. Based on the improved results on the THUMOS’14 dataset, we choose the two-way buffer setting with horizontal flipping of frames for data augmentation.

Results: In Table 3 we show the performance of our model and compare with existing published approaches. The results are shown for two different experimental settings. In the first setting, only the training set is used for training and the performance is shown for either the validation or test data or both. In the second setting, training is performed on both training and validation sets while the performance is shown on the test set. The table shows that the proposed method does achieve a performance better than methods not using handcrafted features *e.g.*, UPC [18]. UPC is the most fair comparison as it also uses only C3D features. However, it relies on a strong assumption that each video in ActivityNet just contains one activity class. Our approach obtains an improvement of 4.3% on the validation set and 4.5% on the test set over UPC [18] in terms of mAP@0.5 without any such strong assumptions. When both training and validation sets are used for training, the performance improves further by 1.6%.

R-C3D falls slightly behind [29] which uses LSTM based tracking and performs activity prediction using deep

Table 4. Temporal Action Localization result on ActivityNet Challenge 2017 in terms of Average-mAP (in percentage).

	train data	validation	test
R-C3D (ours)	train	12.7	13.1
R-C3D (ours)	train+val	-	16.7

features as well as optical flow features from the tracked trajectories. The approach in [30] also uses handcrafted motion features like MBH on top of inception and C3D features in addition to dynamic programing based post processing. However, the heavy use of an ensemble of hand-engineered features and dataset dependent heuristics not only stops these methods from learning in an end-to-end fashion but makes them less general across datasets. Unlike these methods, R-C3D is trainable completely end-to-end and is easily extensible to other datasets with little or no parameter tuning, providing better generalization performance. Our method is also capable of using hand engineered features with a possible boost to performance, and we keep this as a future task. Figure 3(b) shows some representative qualitative results from this dataset.

ActivityNet Challenge 2017: We took part in the ActivityNet Challenge 2017 (Task 4: Temporal Action Localization) under the user name “Huijuan Xu” and the affiliation “Boston University”. We initialized the 3-D convnet of the R-C3D with the model pretrained on Sports-1M and finetuned for ActivityNet classification task. Table 4 shows our results in the new evaluation metric Average-mAP at 10 evenly distributed thresholds between 0.5 and 0.95. We combine the training and validation sets to train our model and get a final mAP result of 16.7% in the localization task.

4.3. Experiments on Charades

Charades [26] is a recently introduced dataset for activity classification and detection tasks. The activity detection task involves daily life activities from 157 classes. The dataset consists of 7985 train and 1863 test videos. The videos are recorded by Amazon Mechanical Turk users based on provided scripts. Apart from low illumination, diversity and casual nature of the videos containing day-to-day activities, an additional challenge of this dataset is the abundance of overlapping activities, sometimes multiple activities having exactly the same start and end times (typical examples include pairs of activities like ‘holding a phone’ and ‘playing with a phone’ or ‘holding a towel’ and ‘tidying up a towel’).

Experimental Setup: For this dataset we sample frames at 5 fps, and the input buffer is set to contain 768 frames. This makes the duration of the buffer approximately 154 seconds which covers all the ground truth activity segments in Charades train set. Keeping the longer duration of the activity segments in mind, for this dataset we choose the number of anchor segments K to be 18 with specific scale values [1,

Table 5. Activity detection results on Charades (in percentage). We report the results using the same evaluation metric as in [25].

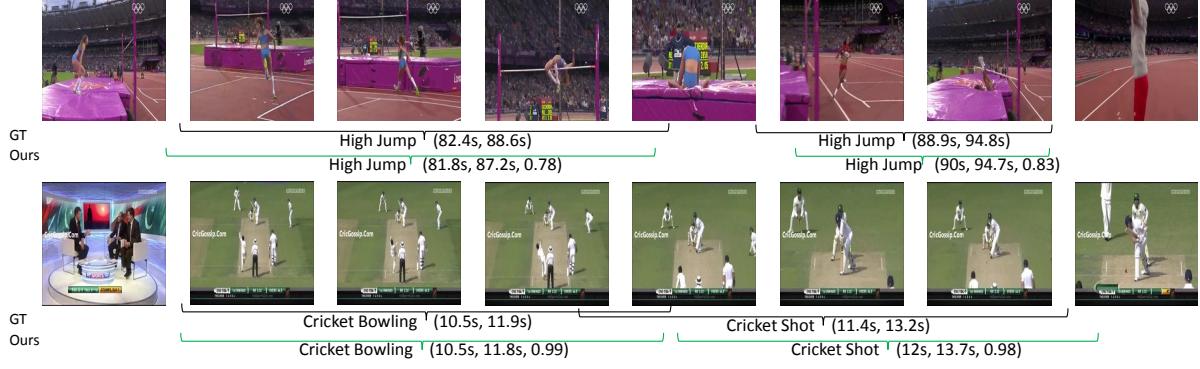
	mAP	
	standard	post-process
Random [25]	4.2	4.2
RGB [25]	7.7	8.8
Two-Stream [25]	7.7	10.0
Two-Stream+LSTM [25]	8.3	8.8
Sigurdsson et al. [25]	9.6	12.1
R-C3D (ours)	12.4	12.7

2, 3, 4, 5, 6, 7, 8, 10, 12, 14, 16, 20, 24, 28, 32, 40, 48]. So the shortest anchor segment has a duration of 1.6 seconds and the longest anchor segment has a duration of 76.8 seconds. Over 99.96% of the activities in the training set is under 76.8 seconds. For this dataset we explored slightly different settings of the anchor segment scales, but found that our model is not very sensitive to this hyperparameter.

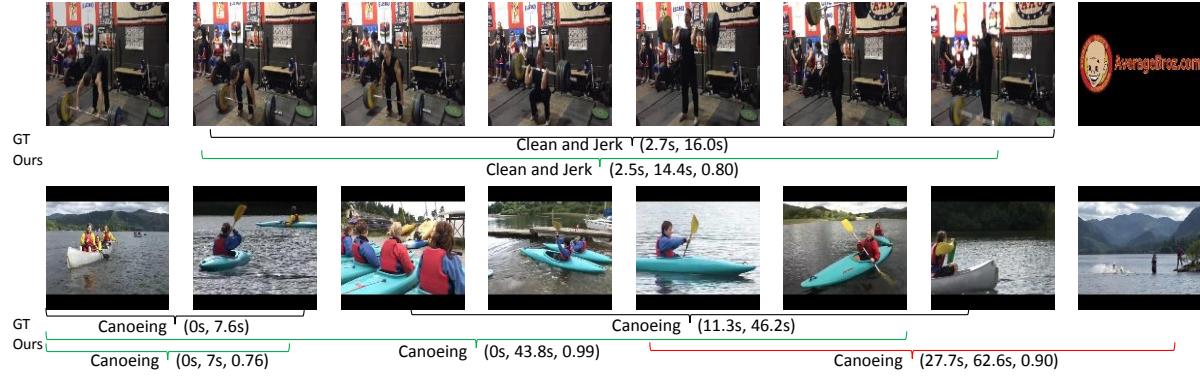
We first finetune the Sports-1M pretrained C3D model [32] on the Charades training set at the same 5 fps and initialize the 3D ConvNet part of our model with these finetuned weights. Next, we train R-C3D end-to-end on Charades by freezing the first two convolutional layers in order to accelerate training. The learning rate is kept fixed at 0.0001 for the first 10 epochs and then decreased to 0.00001 for 5 further epochs. We augment the data by following the two-way buffer setting and horizontal flipping of frames.

Results: Table 5 provides a comparative evaluation of the proposed model with various baseline models reported in [25]. This approach [25] trains a CRF based video classification model (asynchronous temporal fields) and evaluates the prediction performance on 25 equidistant frames by making a multi-label prediction for each of these frames. The activity localization result is reported in terms of mAP metric on these frames. For a fair comparison, we map our activity segment prediction to 25 equidistant frames and evaluate using the same mAP evaluation metric. A second evaluation strategy proposed in this work relies on a post-processing stage where the frame level predictions are averaged across 20 frames leading to more spatial consistency. As shown in the Table 5, our model outperforms the asynchronous temporal fields model proposed in [25] as well as the different baselines reported in the same paper. While the improvement over the standard method is as high as 2.8%, the improvement after the post-processing is not as high. One possible reason could be that our end-to-end fully convolutional model captures the spatial consistency implicitly without requiring any manually-designed postprocessing.

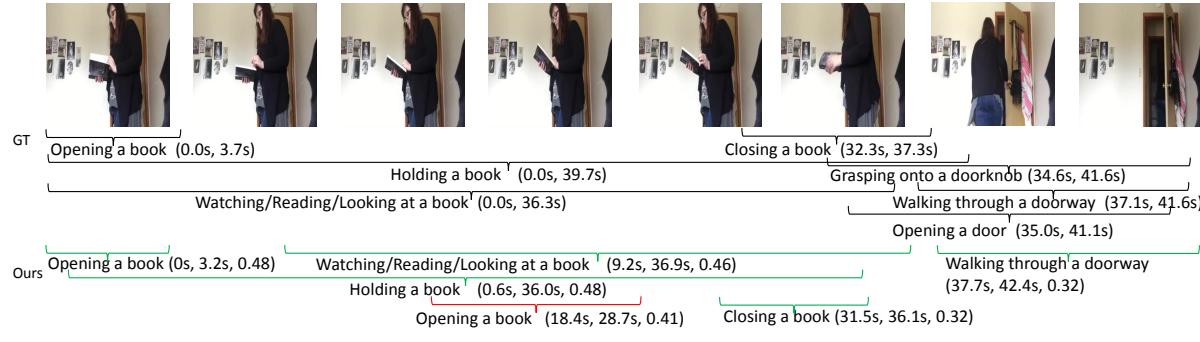
Following the standard practice we also evaluated our model in terms of mAP@0.5 which results in 9.3%. The performance is not at par with other datasets presumably because of the inherent challenges involved in Charades e.g., the low illumination indoor scenes or the multi-label nature of the data. Initialization with a better C3D classification model trained on indoor videos with these chal-



(a) THUMOS'14



(b) ActivityNet



(c) Charades

Figure 3. Qualitative visualization of the predicted activities by R-C3D (best viewed in color). Figure (a) and (b) show results for two videos each in THUMOS'14 and ActivityNet. (c) shows the result for one video from Charades. Groundtruth activity segments are marked in black. Predicted activity segments are marked in green for correct predictions and in red for wrong ones. Predicted activities with IoU more than 0.5 are considered as correct. Corresponding start-end times and confidence score are shown inside brackets.

lenging conditions may further boost the performance. Figure 3(c) shows some representative qualitative results from one video in this dataset.

One of the major challenges of this dataset is the presence of a large number of temporally overlapping activities. The results show that our model is capable of handling such scenarios. This is achieved by the ability of the proposal subnet to produce possibly overlapping activity proposals and is further facilitated by region offset regression.

Table 6. Activity detection speed at test time.

	FPS
S-CNN [24]	60
DAP [4]	134.1
R-C3D (ours on Titan X Maxwell)	569
R-C3D (ours on Titan X Pascal)	1030

4.4. Activity Detection Speed

In this section, we compare our model with two others in terms of detection speed, as shown in Table 6 results.

S-CNN [24] uses a time-consuming sliding window strategy and predicts at 60 fps. DAP [4] incorporates a proposal prediction step on top of LSTM and predicts at 134.1 fps. Our R-C3D model constructs the proposal and classification pipeline in an end-to-end fashion, which is significantly faster at 569 fps on a Titan-X GPU Maxwell. On the upgraded Titan-X GPU Pascal, our test speed reaches an even higher 1030 fps. The speedup of R-C3D over DAP may come from the fact that the LSTM recurrent architecture in DAP takes time to unroll, while R-C3D directly accepts a wide range of frames as input and operates on the shared features in both the proposal and classification subnets.

5. Conclusion

We introduce R-C3D, the first end-to-end temporal proposal classification network for activity detection in untrimmed videos. We evaluate our approach on three large-scale data sets with very diverse characteristics, and demonstrate that it can detect activities faster and more accurately than existing models based on 3D Convnets. Additional features can be incorporated into our model to further boost the activity detection result. One future direction may be to integrate R-C3D with hand-engineered motion features for improved activity prediction without sacrificing the speed.

References

- [1] ActivitNet Large Scale Activity Recognition Challenge. http://activity-net.org/challenges/2016/data/anet_challenge_summary.pdf, 2016. 6
- [2] J. Carreira and C. Sminchisescu. CPMC: Automatic Object Segmentation using Constrained Parametric Min-cuts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7):1312–1328, 2012. 2
- [3] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: Object Detection via Region-based Fully Convolutional Networks. In *Neural Information Processing Systems*, pages 379–387, 2016. 2
- [4] V. Escorcia, F. C. Heilbron, J. C. Niebles, and B. Ghanem. DAPs: Deep Action Proposals for Action Understanding. In *European Conference on Computer Vision*, pages 768–784, 2016. 2, 5, 8, 9
- [5] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object Detection with Discriminatively Trained Part Based Models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32:1627–1645, 2010. 4
- [6] R. Girshick. Fast R-CNN. In *IEEE International Conference on Computer Vision*, pages 1440–1448, 2015. 2, 4
- [7] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 580–587, 2014. 2
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 1
- [9] F. C. Heilbron, V. Escorcia, B. Ghanem, and J. C. Niebles. ActivityNet: A Large-Scale Video Benchmark for Human Activity Understanding. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015. 2, 5, 6
- [10] F. C. Heilbron, J. C. Niebles, and B. Ghanem. Fast Temporal Activity Proposals for Efficient Detection of Human Actions in Untrimmed Videos. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1914–1923, 2016. 1, 5
- [11] S. Ji, W. Xu, M. Yang, and K. Yu. 3D Convolutional Neural Networks for Human Action Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35:221–231, 2013. 2
- [12] Y.-G. Jiang, J. Liu, A. R. Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS Challenge: Action Recognition with a Large Number of Classes. <http://crcv.ucf.edu/THUMOS14/>, 2014. 2, 5
- [13] S. Karaman, L. Seidenari, and A. D. Bimbo. Fast Saliency Based Pooling of Fisher Encoded Dense Trajectories. In *ECCV THUMOS Workshop*, 2014. 1, 2, 5
- [14] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale Video Classification with Convolutional Neural Networks. In *IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014. 2
- [15] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning Realistic Human Actions from Movies. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008. 2
- [16] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single Shot MultiBox Detector. In *European Conference on Computer Vision*, pages 21–37, 2016. 2
- [17] S. Ma, L. Sigal, and S. Sclaroff. Learning Activity Progression in LSTMs for Activity Detection and Early Detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1942–1950, 2016. 2
- [18] A. Montes, A. Salvador, and X. G. i Nieto. Temporal Activity Detection in Untrimmed Videos with Recurrent Neural Networks. *arXiv preprint arXiv:1608.08128*, 2016. 2, 6
- [19] J. Y.-H. Ng, M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici. Beyond Short Snippets: Deep Networks for Video Classification. In *IEEE conference on computer vision and pattern recognition*, pages 4694–4702, 2015. 2
- [20] D. Oneata, J. Verbeek, and C. Schmid. The LEAR submission at Thumos 2014. *ECCV THUMOS Workshop*, 2014. 1, 2, 5
- [21] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *Neural Information Processing Systems*, pages 91–99, 2015. 1, 2, 3, 4
- [22] A. Richard and J. Gall. Temporal Action Detection Using a Statistical Language Model. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2, 5
- [23] Z. Shou, J. Chan, A. Zareian, K. Miyazaway, and S.-F. Chang. CDC: Convolutional-De-Convolutional Networks

- for Precise Temporal Action Localization in Untrimmed Videos. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 5
- [24] Z. Shou, D. Wang, and S.-F. Chang. Temporal Action Localization in Untrimmed Videos via Multi-stage CNNs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 1, 2, 4, 5, 6, 8, 9
- [25] G. A. Sigurdsson, S. Divvala, A. Farhadi, and A. Gupta. Asynchronous Temporal Fields for Action Recognition. *arXiv preprint arXiv:1612.06371*, 2017. 7
- [26] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta. Hollywood in Homes: Crowdsourcing Data Collection for Activity Understanding. In *European Conference on Computer Vision*, 2016. 2, 5, 7
- [27] K. Simonyan and A. Zisserman. Two-stream Convolutional Networks for Action Recognition in Videos. In *Neural Information Processing Systems*, pages 568–576, 2014. 2
- [28] K. Simonyan and A. Zisserman. Very Deep Convolutional Networks for Large-scale Image Recognition. In *IEEE Conference on Learning Representations*, 2015. 1, 2
- [29] B. Singh, T. K. Marks, M. Jones, O. Tuzel, and M. Shao. A Multi-Stream Bi-Directional Recurrent Neural Network for Fine-Grained Action Detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2, 6
- [30] G. Singh and F. Cuzzolin. Untrimmed Video Classification for Activity Detection: submission to ActivityNet Challenge. *arXiv preprint arXiv:1607.01979*, 2016. 6, 7
- [31] C. Szegedy, A. Toshev, and D. Erhan. Deep Neural Networks for Object Detection. In *Neural Information Processing Systems*, pages 2553–2561, 2013. 2
- [32] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning Spatiotemporal Features with 3D Convolutional Networks. In *IEEE International Conference on Computer Vision*, pages 4489–4497, 2015. 1, 2, 3, 5, 6, 7
- [33] H. Wang and C. Schmid. Action Recognition with Improved Trajectories. In *IEEE International Conference on Computer Vision*, pages 3551–3558, 2013. 1, 2
- [34] L. Wang, Y. Qiao, and X. Tang. Video Action Detection with Relational Dynamic-Poselets. In *European Conference on Computer Vision*, pages 565–580, 2014. 1
- [35] L. Wang, Y. Qiao, X. Tang, and L. V. Gool. Actionness Estimation using Hybrid Fully Convolutional Networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2708–2717, 2016. 1
- [36] L. Wang, Y. Xiong, D. Lin, and L. V. Gool. UntrimmedNets for Weakly Supervised Action Recognition and Detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [37] L. Wang, Y. Yu Qiao, and X. Tang. Action Recognition and Detection by Combining Motion and Appearance Features. *ECCV THUMOS Workshop*, 1, 2014. 1, 2, 5
- [38] P. Weinzaepfel, Z. Harchaoui, and C. Schmid. Learning to Track for Spatio-Temporal Action Localization. In *IEEE International Conference on Computer Vision*, 2015. 2
- [39] S. Yeung, O. Russakovsky, G. Mori, and L. Fei-Fei. End-to-end Learning of Action Detection from Frame Glimpses in Videos. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2678–2687, 2016. 2, 4, 5
- [40] G. Yu and J. Yuan. Fast Action Proposals for Human Action Detection and Search. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 2
- [41] J. Yuan, B. Ni, X. Yang, and A. A. Kassim. Temporal Action Localization with Pyramid of Score Distribution Features. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3093–3102, 2016. 2, 5
- [42] J. Zheng, Z. Jiang, and R. Chellappa. Cross-view Action Recognition via Transferable Dictionary Learning. *IEEE Transactions on Image Processing*, 25(6):2542–2556, 2016. 2

Improve Accurate Action Localization by Dense Pose Estimation

Yuxiang Zhou and Jiankang Deng

Imperial College London, UK

{yuxiang.zhou10, j.deng16}@imperial.ac.uk

Abstract

In this paper, we presented a brief description of the proposed method for Activity Detection Task for ActivityNet Challenge 2017. Our method introduces additional classification signals based on the structured segment networks (SSN) [19] and further improved the performance. To be specific, semantic human body segmentation and pose landmarks localization signals are involved in detection progress. Our work in [20] shown that shape-based supervision signals substantially accelerate learning speed, while also improving localization accuracy. We extend dense supervision signals from [20] to apply to all frames of videos alongside with output from SSN to further improve detection accuracy, especially for pose related and sparsely annotated videos as described in fig 4. The method in general achieves state-of-the-art performance on test set and witnesses remarkable improvement on pose related and sparsely annotated categories e.g. sports.

1. Introduction

Activity Detection and temporal action localization [19, 16, 8, 10, 5, 2, 9] has drawn increasing attention to the research community in past few years. Human activity understanding in untrimmed and long videos, especially, are crucial part of real-word applications including video recommendation, video surveillance, human-machine interaction and many others. It is of importance for algorithms to determine not only actions contained in videos but also temporal boundaries (activity starting/ending frames).

However, many methods are trained on short video clips where actions are tightly cropped, while, in practical, videos tend to be long and untrimmed. Driven by ActivityNet [2], a large-scale video benchmark for human activity understanding is released to the research community and consist of 200 activity categories, in which each contains 100 videos collected “in-the-wild”. This dataset brought notable challenges to existing state-of-the-art approaches.

In this paper, we introduce additional classification signals from dense human poses [20] with spatial attentions

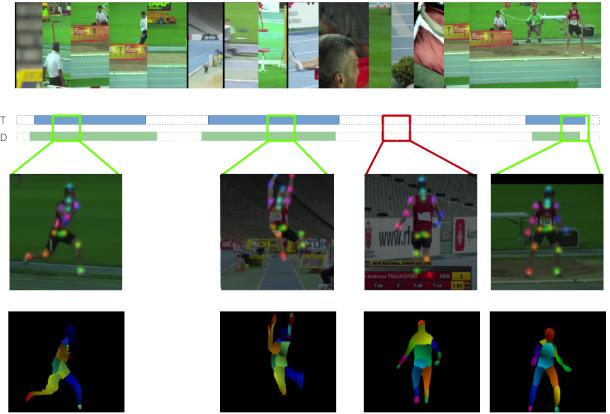


Figure 1. Frames in target video are extracted and applied human detector to generate human proposals before applying our work in [20] to generate heat maps of landmarks and densely correspondent features. These additional signals are combined with the pyramid of action frame proposals from SSN to further improve the performance of determine starting/ending frames. Regarding timelines, **TOP** (blue) shows groundtruth action frame; **BOT** (green) shows predicted action frame. Figure best viewed by zooming in.

alongside with the usage of structured segment networks (SSN) [19]. Specifically, human body detection, semantic human body segmentation and pose landmarks localization signals are involved in activity detection progress. The work in [20] shown that shape-based supervision signals substantially accelerate learning speed, while also improving localization accuracy. We extend dense supervision signals from our work [20] to apply to all frames of videos to generate dense human pose features and combined with results from SSN to further improve accuracy, especially for pose related and sparsely annotated videos as described in Fig 4. The method in general achieves state-of-the-art performance on test set and observed remarkable improvement on pose related and sparsely annotated categories.

2. Structured Segment Networks (SSN)

The SSN network [19] relies on a proposal method (described in section 4) to produce a set of temporal proposals of varying durations, where each proposal comes with a starting and an ending time. Given an input video, a temporal pyramid will be constructed upon each temporal proposals. One proposal is divided into three stages namely *starting*, *course*, and *ending*. In addition to the *course* stage, another level of pyramid with two sub-parts is constructed. To form the global region representations, features from CNNs are pooled within these five parts and concatenated together. The activity classifier and the completeness classifier operate on the the region representations to produce activity probability and class conditional completeness probability. The final probability of the proposal being positive instance is decided by the joint probability from these two classifiers

3. Dense Body Pose Estimation Networks (DenseReg Pose)

Our work in dense body pose estimation networks are inspired by learning with ‘Privileged Information’ [14, 1, 3, 20], where it is argued that one can simplify training through the use of an ‘Intelligent Teacher’ that in a way explains the supervision signal, rather than simply penalizing misclassifications. This technique was recently used in deep learning for the task of image classification [3]; DenseReg Pose shown that shape-based representations provide an excellent source of privileged information for human pose estimation. This additional information is only available during training, only serves as a means of simplifying the training problem, and only requires landmark-level annotations, as all current methods do. Another way of stating this is that we use shape-based representations to construct a set of auxiliary tasks that accelerate and improve the training of pose estimation networks.

Two additional dense supervision signals used in Action Detection task are Support Vector Shape (SVS) [13] and IUV representations [4].

3.1. Support Vector Shape (SVS)

A Support Vector Shape (SVS) is a decision function trained on binary shapes using Support Vector Machines (SVMs) with Radial Basis Function (RBF) kernels [7] - a shape is represented in terms of the classifier’s response on the plane. This representation can be applied to both sparse landmark points and curves, fuses inconsistent landmarks into consistent and directly comparable decision functions, and is robust against noise, missing data, self-occlusions and outliers.

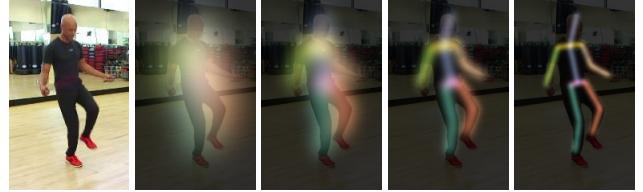


Figure 2. Multichannel Support Vector Shape representations using different granularities. From left to right we show the SVS for $C = [3, 6, 12, 24]$ respectively, where C is the scaling of the underlying SVM data term.

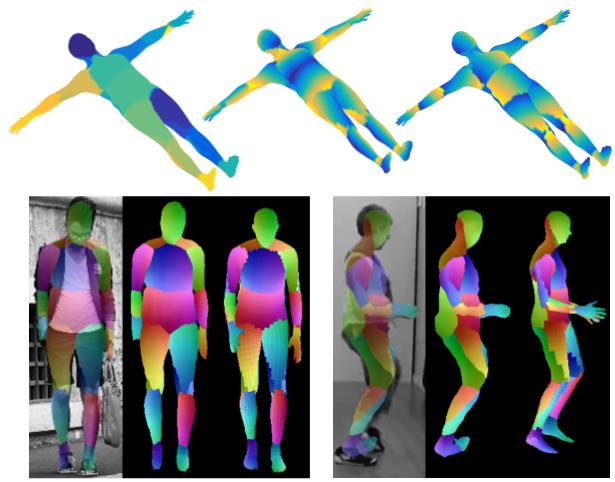


Figure 3. *Top:* Index, U and V fields displayed on the SMPL model. *Bottom:* Dense correspondence results presented as input image fused with estimated UV coordinates, estimated UV coordinates and ground-truth UV coordinates respectively. A customized colour-coding is used for a clear demonstration of correspondence.

3.2. $I - U - V$ representation

We follow [4], who demonstrate that template-to-pixel correspondences can be regressed using a quantized scheme, where each pixel - vertex pair is labelled with an index, I , that indicates membership of a specific patch and a $U - V$ representation as coordinates in a normalized 2D coordinate system. With such a representation, regressing the index I becomes a pixel-wise classification problem and estimating the $U - V$ coordinates becomes an easier pixel-wise continuous regression problem. The challenge in obtaining such a representation on the human shape is the charting: tessellating the shape into patches that are easy to unwrap. To allow an easy unwrapping, we have manually partitioned the surface into 25 patches each of which is isomorphic to the plane. The unwrapping is done using multidimensional scaling (MDS) for each patch, and followed by a normalization to obtain fields $U, V \in [0, 1]$. The I, U and V fields on the SMPL model [6] is presented in Fig. 3

along with CNN-estimated and ground-truth fields for two images that contain frontal and side poses.

4. Activity Detection

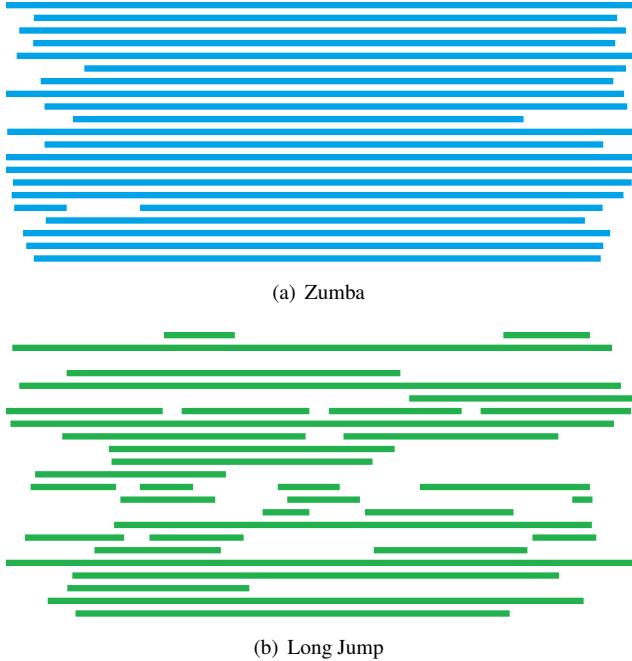


Figure 4. Action annotations of two particular classes on the validation set. Each segment indicates the action duration, which is normalized by the whole video. Detection of the first action (Zumba) is much easier than the second action (Long Jump).

The task for activity detection is to localize the temporal boundary of an activity. There are two types of action annotations as shown in Fig 4, a) action duration is long and almost one action per video in the category e.g. Zumba; b) action duration is short and one video contains multiple actions e.g. Long Jump. Most algorithms tackles the former situation well as these actions are covering majority part of the videos and their activity boundaries are very close to the starting/ending of videos. However, it is challenging to get sufficient accuracy for the second situation where most modern methods failed to provide accurate proposal.

Our methods are focused on determine accurate short activity boundary by incorporating attentions from human detector, poses landmark localization and pose segmentation with activity and completeness classifiers from SSN.

4.1. Temporal Region Proposal

An input video is divided to 20 snippets and temporal region are generated based on sliding windows [9, 17]. A sliding window of size 3 are selected so 18 region proposals are generated. As we incorporated human detector [18], pose estimation and semantic segmentation [20] as

additional feature, existing proposed regions will be duplicated if the human pose detector returns multiple entries. For each proposed region, K-level temporal pyramid where each level dividing the region into smaller parts.

4.2. Temporal Region Classifiers

Structured temporal pyramid pooling [19] is performed to extract global features in which our detection, pose estimation and segmentation are involved. The training and testing of the classifiers are following the SSN network in similar manner.

Two types of linear classifiers (activity classifier and completeness classifier) are implemented on top of high-level features. Given a proposal, the activity classifier will produce a vector of normalized responses via a softmax layer which represents conditional distribution $P(c_i|p_i)$ where c_i is class label and p_i represents given proposal. The completeness classifier C_k are trained for every activity class k . It can be expressed as $P(b_i|c_i, p_i)$ where b_i is binary indicator of the completeness of given region p_i . Outputs of both formed a joint distribution so the loss function is defined as:

$$Loss = -\log P(c_i|p_i) - 1_{(c_i>1)}P(b_i|c_i, p_i)$$

where completeness term is active only when class label is not background.

5. Experiments

The results of our methods on ActivityNet 2017 (test set) are shown in Table 1. Highest ranked results submitted to previous ActivityNet 2016 Challenge and current ActivityNet 2017 Challenge are involved in the table.

ActivityNet 2017 (Test Set)	
Method	Avg. mAP
Wang, R. and Tao, D. [15]	14.62
Singh, B. and Marks, T. et. al. [11]	16.68
Singh, G. and Cuzzolin, F. [12]	17.83
Zhao, Y. and Xiong, Y. et. al. [19]	28.28
Xiong, Y. et. al. [ActivityNet 2017]	31.863
Lin, T. et. al. [ActivityNet 2017]	33.406
Our Method	31.826

Table 1. Action detection results on ActivityNet 2017. Evaluated by mean average precision (mAP).

References

- [1] Unifying distillation and privileged information. In *ICLR*, 2016.
- [2] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970, 2015.

- [3] Y. Chen, X. Jin, J. Feng, and S. Yan. Training group orthogonal neural networks with privileged information. *CoRR*, abs/1701.06772, 2017.
- [4] R. A. Güler, G. Trigeorgis, E. Antonakos, P. Snape, S. Zafeiriou, and I. Kokkinos. Densereg: Fully convolutional dense shape regression in-the-wild. *arXiv preprint arXiv:1612.01202*, 2016.
- [5] Y. Jiang, J. Liu, A. R. Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. Thumos challenge: Action recognition with a large number of classes, 2014.
- [6] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)*, 34(6):248, 2015.
- [7] H. V. Nguyen and F. Porikli. Support vector shape: A classifier-based shape representation. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2013.
- [8] Z. Shou, J. Chan, A. Zareian, K. Miyazawa, and S.-F. Chang. Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. *arXiv preprint arXiv:1703.01515*, 2017.
- [9] Z. Shou, D. Wang, and S.-F. Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1049–1058, 2016.
- [10] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*, pages 568–576, 2014.
- [11] B. Singh, T. K. Marks, M. Jones, O. Tuzel, and M. Shao. A multi-stream bi-directional recurrent neural network for fine-grained action detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1961–1970, 2016.
- [12] G. Singh and F. Cuzzolin. Untrimmed video classification for activity detection: submission to activitynet challenge. *arXiv preprint arXiv:1607.01979*, 2016.
- [13] H. Van Nguyen and F. Porikli. Support vector shape: A classifier-based shape representation. *IEEE transactions on pattern analysis and machine intelligence*, 35(4):970–982, 2013.
- [14] V. Vapnik and A. Vashist. A new learning paradigm: Learning using privileged information. *Neural Networks*, 2009.
- [15] R. Wang and D. Tao. Uts at activitynet 2016. 2016.
- [16] H. Xu, A. Das, and K. Saenko. R-c3d: Region convolutional 3d network for temporal activity detection. *arXiv preprint arXiv:1703.07814*, 2017.
- [17] J. Yuan, B. Ni, X. Yang, and A. A. Kassim. Temporal action localization with pyramid of score distribution features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3093–3102, 2016.
- [18] X. Zeng, W. Ouyang, J. Yan, H. Li, T. Xiao, K. Wang, Y. Liu, Y. Zhou, B. Yang, Z. Wang, et al. Crafting gbd-net for object detection. *arXiv preprint arXiv:1610.02579*, 2016.
- [19] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, D. Lin, and X. Tang. Temporal action detection with structured segment networks. *arXiv preprint arXiv:1704.06228*, 2017.
- [20] Y. Zhou, A. Guler, G. Trigeorgis, I. Kokkinos, and S. Zafeiriou. Training body pose estimation networks with dense shape supervision. 2017.

RUC-CMU: System Descriptions for the Dense Video Captioning Task

Qin Jin

Renmin University of China

qjin@ruc.edu.cn

Shizhe Chen

Renmin University of China

cszhe1@ruc.edu.cn

Jia Chen

Carnegie Mellon University

jiac@cs.cmu.edu

Alexander Hauptmann

Carnegie Mellon University

alex@cs.cmu.edu

1. Introduction

In this paper, we present our approaches in the ActivityNet Dense Video Captioning Task. We explore different caption models including the mean pooling model, temporal attention model and incorporate contexts into the caption models. Both the temporal attention mechanism and contexts have improved the caption performance on the validation set, and our temporal attention model with contexts achieves 9.62 METEOR score in the challenge testing set.

2. Approaches

2.1. Features

Besides the motion features C3D provided by the challenge organizers, we extract static image features from the last pooling layer of the resnet200 pretrained on imagenet11k dataset [1] every 8 frames.

2.2. Caption Models

Mean Pooling Caption Model (MP) [3] is the baseline model for video captioning task, which consists of a multimodal fusion encoder and vanilla LSTM decoder. The multimodal encoder applies mean pooling over the sequential image and motion features, and then uses a neural network to map the feature concatenations into a dense video representation x . We set x as the initial state h_0 of the LSTM decoder and then generate word sequences as follows:

$$h_t = f(h_{t-1}, w_{t-1}) \text{ for } t = 1, \dots, N_w \quad (1)$$

where f is the LSTM update function [2], h_t is the state of LSTM and w_0 is the start symbol. Then the probability of the correct word conditioned on the video content and previous words can be expressed as:

$$\Pr(w_t|x, w_0, \dots, w_{t-1}) = \text{Softmax}(W_d h_t + b_d) \quad (2)$$

where W_d, b_d are parameters to be learned.

Temporal Attention Caption Model (TA) [4] can pay attention to relevant temporal information when generating different words. Assume $\{v_1, \dots, v_n\}$ is the set of visual features in different timesteps, where n is the total visual timesteps. When predicting the t -th word, the set of visual features are weighted by different attention scores $\alpha_i^{(t)}$ to obtain the relevant visual context feature $\phi_t(v)$ as follows:

$$e_i^{(t)} = w^T \tanh(W_a h_{t-1} + U_a v_i + b_a) \quad (3)$$

$$\alpha_i^{(t)} = \exp\{e_i^{(t)}\} / \sum_{j=1}^n \exp\{e_j^{(t)}\} \quad (4)$$

$$\phi_t(v) = \sum_{i=1}^N \alpha_i^{(t)} v_i \quad (5)$$

where w, W_a, U_a and b_a are the parameters to be learned. $\phi_t(v)$ is then concatenated with the previous word embedding as the input of the LSTM decoder:

$$h_t = f(h_{t-1}, [\cdot; \phi_t(v)]; \theta_d) \text{ for } t = 1, \dots, N_w \quad (6)$$

where $[\cdot; \cdot]$ denotes the concatenation operation and the h_0 and w_0 are set in the same way as the MP model.

Temporal Attention Caption Model with Contexts (TAC): Since the description for the current event is related to its neighboring events, we adopt a simple strategy to incorporate the neighboring event contexts into the caption model. We apply mean pooling over the features that are extracted from frames before and after the reference event respectively as the context features. Then we concatenate the context features together with features within the current event to initialize the hidden state of LSTM language decoder. The temporal features to be attended on are still the features within the current event. In this way, the caption model can be aware of the event context but still focus on the current event when generating the event descriptions.

We compare the different models using the groundtruth proposal in the validation set. As shown in the Table 1, the

Table 1. Caption performance using different caption models.

	BLEU@4	METEOR	ROUGHL	CIDEr
MP	2.24	10.50	20.55	28.41
TA	2.38	10.79	21.13	28.93
TAC	2.53	10.67	21.74	31.39

Table 2. Performance of the even temporal proposals on the validation set.

tIoU	0.3	0.5	0.7	Average
Recall	97.71	94.02	77.91	89.88
Precision	91.42	52.41	24.19	56.01

Table 3. Caption performance comparison of the TAC model using groundtruth and even temporal proposals.

	BLEU@4	METEOR	ROUGHL	CIDEr
groundtruth	2.53	10.67	21.74	31.39
even	2.70	10.62	20.60	27.70

TA model outperforms the MP model on all four metrics, and the context features are beneficial which improve the caption performance of the TA model on all the metrics except slight decrease on the METEOR score.

2.3. Temporal Proposals

As a simple baseline, we divide the video into even segments as the temporal proposals. To generate proposals in different temporal scales, we utilize five different levels for segmentation which results in 15 temporal proposals for each video. Table 2 presents the performance of the our temporal proposals on the validation set. We can see

that temporal proposals generated from the simple division method can achieve high coverage of the groundtruth proposals with average 89.88% recall. We also explore some off-the-shelf video segmentation toolkits, but all of them are inferior to the even temporal proposals. This might because the off-the-shelf toolkits segment the videos according to scene changes but not the semantic event changes.

We also compare the caption performance using the groundtruth proposals and the even temporal proposals in Table 3, which also suggests the even segmentation is a strong baseline for the task.

2.4. Submission Results

We submit the captions predicted by our TAC model with even temporal proposals. The METEOR score on the testing set of the challenge is 9.62.

References

- [1] Resnet200 image features. <http://data.dmlc.ml/mxnet/models/>, 2017.
- [2] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [3] Q. Jin, J. Chen, S. Chen, Y. Xiong, and A. Hauptmann. Describing videos using multi-modal fusion. In *ACM*, pages 1087–1091, 2016.
- [4] L. Yao, A. Torabi, K. Cho, N. Ballas, C. J. Pal, H. Larochelle, and A. C. Courville. Describing videos by exploiting temporal structure. In *ICCV*, pages 4507–4515, 2015.